
d3rlpy

Takuma Seno

Jul 20, 2021

TUTORIALS

1	Getting Started	3
1.1	Install	3
1.2	Prepare Dataset	3
1.3	Setup Algorithm	4
1.4	Setup Metrics	4
1.5	Start Training	5
1.6	Save and Load	6
2	Jupyter Notebooks	7
3	API Reference	9
3.1	Algorithms	9
3.2	Q Functions	245
3.3	MDPDataSet	251
3.4	Datasets	262
3.5	Preprocessing	265
3.6	Optimizers	273
3.7	Network Architectures	276
3.8	Metrics	283
3.9	Off-Policy Evaluation	291
3.10	Save and Load	311
3.11	Logging	313
3.12	scikit-learn compatibility	314
3.13	Online Training	316
3.14	Model-based Algorithms	323
3.15	Stable-Baselines3 Wrapper	331
4	Command Line Interface	333
4.1	plot	333
4.2	plot-all	334
4.3	export	335
4.4	record	335
4.5	play	336
5	Installation	337
5.1	Recommended Platforms	337
5.2	Install d3rlpy	337
6	Tips	339
6.1	Reproducibility	339
6.2	Learning from image observation	339

6.3	Improve performance beyond the original paper	340
7	Paper Reproductions	341
7.1	Offline	341
7.2	Online	343
8	License	345
9	Indices and tables	347
	Python Module Index	349
	Index	351

d3rlpy is a easy-to-use offline deep reinforcement learning library.

```
$ pip install d3rlpy
```

d3rlpy provides state-of-the-art offline deep reinforcement learning algorithms through out-of-the-box scikit-learn-style APIs. Unlike other RL libraries, the provided algorithms can achieve extremely powerful performance beyond their papers via several tweaks.

GETTING STARTED

This tutorial is also available on [Google Colaboratory](#)

1.1 Install

First of all, let's install `d3rlpy` on your machine:

```
$ pip install d3rlpy
```

See more information at [Installation](#).

Note: If `core dump` error occurs in this tutorial, please try [Install from source](#).

Note: `d3rlpy` supports Python 3.6+. Make sure which version you use.

Note: If you use GPU, please setup CUDA first.

1.2 Prepare Dataset

You can make your own dataset without any efforts. In this tutorial, let's use integrated datasets to start. If you want to make a new dataset, see [MDP Dataset](#).

`d3rlpy` provides suites of datasets for testing algorithms and research. See more documents at [Datasets](#).

```
from d3rlpy.datasets import get_cartpole # CartPole-v0 dataset
from d3rlpy.datasets import get_pendulum # Pendulum-v0 dataset
from d3rlpy.datasets import get_pybullet # PyBullet task datasets
from d3rlpy.datasets import get_atari    # Atari 2600 task datasets
from d3rlpy.datasets import get_d4rl     # D4RL datasets
```

Here, we use the CartPole dataset to instantly check training results.

```
dataset, env = get_cartpole()
```

One interesting feature of `d3rlpy` is full compatibility with scikit-learn utilities. You can split dataset into a training dataset and a test dataset just like supervised learning as follows.

```
from sklearn.model_selection import train_test_split

train_episodes, test_episodes = train_test_split(dataset, test_size=0.2)
```

1.3 Setup Algorithm

There are many algorithms available in d3rlpy. Since CartPole is the simple task, let's start from DQN, which is the Q-learning algorithm proposed as the first deep reinforcement learning algorithm.

```
from d3rlpy.algos import DQN

# if you don't use GPU, set use_gpu=False instead.
dqn = DQN(use_gpu=True)

# initialize neural networks with the given observation shape and action size.
# this is not necessary when you directly call fit or fit_online method.
dqn.build_with_dataset(dataset)
```

See more algorithms and configurations at [Algorithms](#).

1.4 Setup Metrics

Collecting evaluation metrics is important to train algorithms properly. In d3rlpy, the metrics is computed through scikit-learn style scorer functions.

```
from d3rlpy.metrics.scorer import td_error_scorer
from d3rlpy.metrics.scorer import average_value_estimation_scorer

# calculate metrics with test dataset
td_error = td_error_scorer(dqn, test_episodes)
```

Since evaluating algorithms without access to environment is still difficult, the algorithm can be directly evaluated with `evaluate_on_environment` function if the environment is available to interact.

```
from d3rlpy.metrics.scorer import evaluate_on_environment

# set environment in scorer function
evaluate_scorer = evaluate_on_environment(env)

# evaluate algorithm on the environment
rewards = evaluate_scorer(dqn)
```

See more metrics and configurations at [Metrics](#).

1.5 Start Training

Now, you have all to start data-driven training.

```
dqn.fit(train_episodes,
        eval_episodes=test_episodes,
        n_epochs=10,
        scorers={
            'td_error': td_error_scorer,
            'value_scale': average_value_estimation_scorer,
            'environment': evaluate_scorer
        })
```

Then, you will see training progress in the console like below:

```
augmentation=[]
batch_size=32
bootstrap=False
dynamics=None
encoder_params={}
eps=0.00015
gamma=0.99
learning_rate=6.25e-05
n_augmentations=1
n_critics=1
n_frames=1
q_func_factory=mean
scaler=None
share_encoder=False
target_update_interval=8000.0
use_batch_norm=True
use_gpu=None
observation_shape=(4,)
action_size=2
100%| 2490/2490 [00:24<00:00, 100.63it/s]
epoch=0 step=2490 value_loss=0.190237
epoch=0 step=2490 td_error=1.483964
epoch=0 step=2490 value_scale=1.241220
epoch=0 step=2490 environment=157.400000
100%| 2490/2490 [00:24<00:00, 100.63it/s]
.
.
.
```

See more about logging at [Logging](#).

Once the training is done, your algorithm is ready to make decisions.

```
observation = env.reset()

# return actions based on the greedy-policy
action = dqn.predict([observation])[0]

# estimate action-values
value = dqn.predict_value([observation], [action])[0]
```

1.6 Save and Load

d3rlpy provides several ways to save trained models.

```
# save full parameters
dqn.save_model('dqn.pt')

# load full parameters
dqn2 = DQN()
dqn2.build_with_dataset(dataset)
dqn2.load_model('dqn.pt')

# save the greedy-policy as TorchScript
dqn.save_policy('policy.pt')

# save the greedy-policy as ONNX
dqn.save_policy('policy.onnx', as_onnx=True)
```

See more information at *Save and Load*.

JUPYTER NOTEBOOKS

- CartPole
- Toy task (line tracer)
- Continuous Control with PyBullet
- Discrete Control with Atari

API REFERENCE

3.1 Algorithms

d3rlpy provides state-of-the-art offline deep reinforcement learning algorithms as well as online algorithms for the base implementations.

3.1.1 Continuous control algorithms

<code>d3rlpy.algos.BC</code>	Behavior Cloning algorithm.
<code>d3rlpy.algos.DDPG</code>	Deep Deterministic Policy Gradients algorithm.
<code>d3rlpy.algos.TD3</code>	Twin Delayed Deep Deterministic Policy Gradients algorithm.
<code>d3rlpy.algos.SAC</code>	Soft Actor-Critic algorithm.
<code>d3rlpy.algos.BCQ</code>	Batch-Constrained Q-learning algorithm.
<code>d3rlpy.algos.BEAR</code>	Bootstrapping Error Accumulation Reduction algorithm.
<code>d3rlpy.algos.CRR</code>	Critic Regularized Regression algorithm.
<code>d3rlpy.algos.CQL</code>	Conservative Q-Learning algorithm.
<code>d3rlpy.algos.AWR</code>	Advantage-Weighted Regression algorithm.
<code>d3rlpy.algos.AWAC</code>	Advantage Weighted Actor-Critic algorithm.
<code>d3rlpy.algos.PLAS</code>	Policy in Latent Action Space algorithm.
<code>d3rlpy.algos.PLASWithPerturbation</code>	Policy in Latent Action Space algorithm with perturbation layer.
<code>d3rlpy.algos.MOPO</code>	Model-based Offline Policy Optimization.
<code>d3rlpy.algos.COMBO</code>	Conservative Offline Model-Based Optimization.
<code>d3rlpy.algos.RandomPolicy</code>	Random Policy for continuous control algorithm.

d3rlpy.algos.BC

```
class d3rlpy.algos.BC (*, learning_rate=0.001, optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), en-
coder_factory='default', batch_size=100, n_frames=1, use_gpu=False,
scaler=None, action_scaler=None, impl=None, **kwargs)
```

Behavior Cloning algorithm.

Behavior Cloning (BC) is to imitate actions in the dataset via a supervised learning approach. Since BC is only imitating action distributions, the performance will be close to the mean of the dataset even though BC mostly

works better than online RL algorithms.

$$L(\theta) = \mathbb{E}_{a_t, s_t \sim D} [(a_t - \pi_\theta(s_t))^2]$$

Parameters

- **learning_rate** (*float*) – learning rate.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **use_gpu** (*bool*, *int* or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessing.Scaler` or *str*) – preprocessor. The available options are [`'pixel'`, `'min_max'`, `'standard'`].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or *str*) – action scaler. The available options are [`'min_max'`].
- **impl** (`d3rlpy.algos.torch.bc_impl.BCImpl`) – implementation of the algorithm.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (*Optional* [`d3rlpy.online.buffers.Buffer`]) – replay buffer.
- **explorer** (*Optional* [`d3rlpy.online.explorers.Explorer`]) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.

- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)
Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.

- **scorers** (*Optional*[*Dict*[*str*, *Callable*[[*Any*, *List*[*d3rlpy.dataset.Episode*]], *float*]]]) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type *List*[*Tuple*[*int*, *Dict*[*str*, *float*]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-  
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-  
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffers.BatchBuffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.

- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fit_online (*env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000, update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If `None`, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If `False`, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*) , which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type `Generator[Tuple[int, Dict[str, float]], None, None]`

classmethod `from_json` (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (fname)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (str) – source file path.

Return type None

predict (x)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (Union[numpy.ndarray, List[Any]]) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (x, action, with_std=False)

value prediction is not supported by BC algorithms.

Parameters

- **x** (Union[numpy.ndarray, List[Any]]) –
- **action** (Union[numpy.ndarray, List[Any]]) –
- **with_std** (bool) –

Return type numpy.ndarray

sample_action (x)

sampling action is not supported by BC algorithm.

Parameters **x** (Union[numpy.ndarray, List[Any]]) –

Return type None

save_model (fname)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (str) – destination file path.

Return type None

save_params (logger)

Saves configurations as params.json.

Parameters `logger` (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DDPG

```
class d3rlpy.algos.DDPG(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003, ac-
    tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    actor_encoder_factory='default', critic_encoder_factory='default',
    q_func_factory='mean', batch_size=100, n_frames=1, n_steps=1,
    gamma=0.99, tau=0.005, n_critics=1, target_reduction_type='min',
    use_gpu=False, scaler=None, action_scaler=None, impl=None,
    **kwargs)
```

Deep Deterministic Policy Gradients algorithm.

DDPG is an actor-critic algorithm that trains a Q function parametrized with θ and a policy function parametrized with ϕ .

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} \left[\left(r_{t+1} + \gamma Q_{\theta'}(s_{t+1}, \pi_{\phi'}(s_{t+1})) - Q_{\theta}(s_t, a_t) \right)^2 \right]$$

$$J(\phi) = \mathbb{E}_{s_t \sim D} \left[Q_{\theta}(s_t, \pi_{\phi}(s_t)) \right]$$

where θ' and ϕ are the target network parameters. There target network parameters are updated every iteration.

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta'$$

$$\phi' \leftarrow \tau \phi + (1 - \tau) \phi'$$

References

- Silver et al., Deterministic policy gradient algorithms.
- Lillicrap et al., Continuous control with deep reinforcement learning.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q function.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the critic.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.

- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.ddpg_impl.DDPGImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated flag is True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffer.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type *None*

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)
Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-  
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-  
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (`d3rlpy.envs.batch.BatchEnv`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.BatchBuffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_epochs** (`int`) – the number of epochs to train.
- **n_steps_per_epoch** (`int`) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (`int`) – the number of updates per epoch.
- **eval_interval** (`int`) – the number of epochs before evaluation.
- **eval_env** (`Optional[gym.core.Env]`) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (`float`) – ϵ -greedy factor during evaluation.
- **save_metrics** (`bool`) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (`int`) – the number of epochs before saving models.
- **experiment_name** (`Optional[str]`) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (`bool`) – flag to add timestamp string to the last of directory name.
- **logdir** (`str`) – root directory name to save logs.
- **verbose** (`bool`) – flag to show logged information on stdout.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **tensorboard_dir** (`Optional[str]`) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (`bool`) – flag to turn terminal flag False when `TimeLimit.truncated` flag is True, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (`Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]`) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type *None*

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.TD3

```
class d3rlpy.algos.TD3(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), actor_encoder_factory='default', critic_encoder_factory='default',
q_func_factory='mean', batch_size=100, n_frames=1, n_steps=1, gamma=0.99, tau=0.005, n_critics=2, target_reduction_type='min',
target_smoothing_sigma=0.2, target_smoothing_clip=0.5, update_actor_interval=2, use_gpu=False, scaler=None, action_scaler=None,
impl=None, **kwargs)
```

Twin Delayed Deep Deterministic Policy Gradients algorithm.

TD3 is an improved DDPG-based algorithm. Major differences from DDPG are as follows.

- TD3 has twin Q functions to reduce overestimation bias at TD learning. The number of Q functions can be designated by `n_critics`.
- TD3 adds noise to target value estimation to avoid overfitting with the deterministic policy.
- TD3 updates the policy function after several Q function updates in order to reduce variance of action-value estimation. The interval of the policy function update can be designated by `update_actor_interval`.

$$L(\theta_i) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} [(r_{t+1} + \gamma \min_j Q_{\theta'_j}(s_{t+1}, \pi_{\phi'}(s_{t+1}) + \epsilon) - Q_{\theta_i}(s_t, a_t))^2]$$

$$J(\phi) = \mathbb{E}_{s_t \sim D} [\min_i Q_{\theta_i}(s_t, \pi_{\phi}(s_t))]$$

where $\epsilon \sim \text{clip}(N(0, \sigma), -c, c)$

References

- Fujimoto et al., Addressing Function Approximation Error in Actor-Critic Methods.

Parameters

- **actor_learning_rate** (`float`) – learning rate for a policy function.
- **critic_learning_rate** (`float`) – learning rate for Q functions.
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.

- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the critic.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **tau** (`float`) – target network synchronization coefficient.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_reduction_type** (`str`) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **target_smoothing_sigma** (`float`) – standard deviation for target noise.
- **target_smoothing_clip** (`float`) – clipping range for target noise.
- **update_actor_interval** (`int`) – interval to update policy function described as *delayed policy update* in the paper.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or `str`) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or `str`) – action preprocessor. The available options are ['min_max'].
- **impl** (`d3rlpy.algos.torch.td3_impl.TD3Impl`) – algorithm implementation.

Methods

build_with_dataset (`dataset`)

Instantiate implementation object with MDPDataset object.

Parameters `dataset` (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (`env`)

Instantiate implementation object with OpenAI Gym object.

Parameters `env` (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (`env`, `buffer=None`, `explorer=None`, `n_steps=1000000`, `show_progress=True`, `time_limit_aware=True`)

Collects data via interaction with environment.

If `buffer` is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (*`gym.core.Env`*) – gym-like environment.
- **buffer** (*`Optional[d3rlpy.online.buffers.Buffer]`*) – replay buffer.
- **explorer** (*`Optional[d3rlpy.online.explorers.Explorer]`*) – action explorer.
- **n_steps** (*`int`*) – the number of total steps to train.
- **show_progress** (*`bool`*) – flag to show progress bar for iterations.
- **timelimit_aware** (*`bool`*) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*`observation_shape, action_size`*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when `fit` method is called.

Parameters

- **observation_shape** (*`Sequence[int]`*) – observation shape.
- **action_size** (*`int`*) – dimension of action-space.

Return type `None`

fit (*`dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None`*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*`Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]`*) – list of episodes to train.
- **n_epochs** (*`Optional[int]`*) – the number of epochs to train.
- **n_steps** (*`Optional[int]`*) – the number of steps to train.
- **n_steps_per_epoch** (*`int`*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*`bool`*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*`Optional[str]`*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*`bool`*) – flag to add timestamp string to the last of directory name.
- **logdir** (*`str`*) – root directory name to save logs.

- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.

- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Return type *None*

fit_online (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *n_steps_per_epoch=10000*, *update_interval=1*, *update_start_step=0*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.

- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type None

fitter (dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when n_steps is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with eval_episodes.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.

- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type *Generator*[*Tuple*[*int*, *Dict*[*str*, *float*]], *None*, *None*]

classmethod from_json (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional*[*Union*[*bool*, *int*, *d3rlpy.gpu.Device*]]) – flag to use GPU, device ID or device.

Returns algorithm.

Return type *d3rlpy.base.LearnableBase*

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List*[*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters `deep` (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type `Dict[str, Any]`

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters `fname` (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters `x` (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.SAC

```
class d3rlpy.algos.SAC(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003,
    temp_learning_rate=0.0003, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    actor_encoder_factory='default', critic_encoder_factory='default',
    q_func_factory='mean', batch_size=100, n_frames=1, n_steps=1,
    gamma=0.99, tau=0.005, n_critics=2, target_reduction_type='min', initial_temperature=1.0, use_gpu=False, scaler=None, action_scaler=None,
    impl=None, **kwargs)
```

Soft Actor-Critic algorithm.

SAC is a DDPG-based maximum entropy RL algorithm, which produces state-of-the-art performance in online RL settings. SAC leverages twin Q functions proposed in TD3. Additionally, *delayed policy update* in TD3 is also implemented, which is not done in the paper.

$$L(\theta_i) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D, a_{t+1} \sim \pi_{\phi}(\cdot | s_{t+1})} \left[(y - Q_{\theta_i}(s_t, a_t))^2 \right]$$

$$y = r_{t+1} + \gamma \left(\min_j Q_{\theta_j}(s_{t+1}, a_{t+1}) - \alpha \log(\pi_{\phi}(a_{t+1} | s_{t+1})) \right)$$

$$J(\phi) = \mathbb{E}_{s_t \sim D, a_t \sim \pi_{\phi}(\cdot | s_t)} \left[\alpha \log(\pi_{\phi}(a_t | s_t)) - \min_i Q_{\theta_i}(s_t, \pi_{\phi}(a_t | s_t)) \right]$$

The temperature parameter α is also automatically adjustable.

$$J(\alpha) = \mathbb{E}_{s_t \sim D, a_t \sim \pi_\phi(\cdot|s_t)} \left[-\alpha \left(\log(\pi_\phi(a_t|s_t)) + H \right) \right]$$

where H is a target entropy, which is defined as $\dim a$.

References

- Haarnoja et al., Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor.
- Haarnoja et al., Soft Actor-Critic Algorithms and Applications.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **temp_learning_rate** (*float*) – learning rate for temperature parameter.
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.
- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **temp_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the temperature.
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the critic.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **initial_temperature** (*float*) – initial temperature value.
- **use_gpu** (*bool*, *int* or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].

- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.sac_impl.SACImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If buffer is not given, ReplayBuffer will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffers.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type *None*

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None,
                  with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True,
                  callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.

- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.

- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(episodes):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.

- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type Optional[List[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[*str*, Any]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List*[Any]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (*Union[numpy.ndarray, List[Any]]*) – observations
- **action** (*Union[numpy.ndarray, List[Any]]*) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type *Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]*

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.BCQ

```

class d3rlpy.algos.BCQ(*, actor_learning_rate=0.001, critic_learning_rate=0.001, imita-
tor_learning_rate=0.001, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), imita-
tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), ac-
tor_encoder_factory='default', critic_encoder_factory='default', imita-
tor_encoder_factory='default', q_func_factory='mean', batch_size=100,
n_frames=1, n_steps=1, gamma=0.99, tau=0.005, n_critics=2,
update_actor_interval=1, lam=0.75, n_action_samples=100, ac-
tion_flexibility=0.05, rl_start_step=0, latent_size=32, beta=0.5,
use_gpu=False, scaler=None, action_scaler=None, impl=None, **kwargs)

```

Batch-Constrained Q-learning algorithm.

BCQ is the very first practical data-driven deep reinforcement learning algorithm. The major difference from DDPG is that the policy function is represented as combination of conditional VAE and perturbation function in order to remedy extrapolation error emerging from target value estimation.

The encoder and the decoder of the conditional VAE is represented as E_ω and D_ω respectively.

$$L(\omega) = E_{s_t, a_t \sim D}[(a - \tilde{a})^2 + D_{KL}(N(\mu, \sigma) | N(0, 1))]$$

where $\mu, \sigma = E_\omega(s_t, a_t)$, $\tilde{a} = D_\omega(s_t, z)$ and $z \sim N(\mu, \sigma)$.

The policy function is represented as a residual function with the VAE and the perturbation function represented as $\xi_\phi(s, a)$.

$$\pi(s, a) = a + \Phi \xi_\phi(s, a)$$

where $a = D_\omega(s, z)$, $z \sim N(0, 0.5)$ and Φ is a perturbation scale designated by *action_flexibility*. Although the policy is learned closely to data distribution, the perturbation function can lead to more rewarded states.

BCQ also leverages twin Q functions and computes weighted average over maximum values and minimum values.

$$L(\theta_i) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D}[(y - Q_{\theta_i}(s_t, a_t))^2]$$

$$y = r_{t+1} + \gamma \max_{a_i} [\lambda \min_j Q_{\theta'_j}(s_{t+1}, a_i) + (1 - \lambda) \max_j Q_{\theta'_j}(s_{t+1}, a_i)]$$

where $\{a_i \sim D(s_{t+1}, z), z \sim N(0, 0.5)\}_{i=1}^n$. The number of sampled actions is designated with *n_action_samples*.

Finally, the perturbation function is trained just like DDPG's policy function.

$$J(\phi) = \mathbb{E}_{s_t \sim D, a_t \sim D_\omega(s_t, z), z \sim N(0, 0.5)}[Q_{\theta_1}(s_t, \pi(s_t, a_t))]$$

At inference time, action candidates are sampled as many as *n_action_samples*, and the action with highest value estimation is taken.

$$\pi'(s) = \operatorname{argmax}_{\pi(s, a_i)} Q_{\theta_1}(s, \pi(s, a_i))$$

Note: The greedy action is not deterministic because the action candidates are always randomly sampled. This might affect *save_policy* method and the performance at production.

References

- Fujimoto et al., Off-Policy Deep Reinforcement Learning without Exploration.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **imitator_learning_rate** (*float*) – learning rate for Conditional VAE.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **imitator_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the conditional VAE.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the critic.
- **imitator_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the conditional VAE.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory or str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **update_actor_interval** (*int*) – interval to update policy function.
- **lam** (*float*) – weight factor for critic ensemble.
- **n_action_samples** (*int*) – the number of action samples to estimate action-values.
- **action_flexibility** (*float*) – output scale of perturbation function represented as Φ .
- **rl_start_step** (*int*) – step to start to update policy function and Q functions. If this is large, RL training would be more stabilized.
- **latent_size** (*int*) – size of latent vector for Conditional VAE.
- **beta** (*float*) – KL regularization term for Conditional VAE.
- **use_gpu** (*bool, int or d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler or str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].

- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.bcq_impl.BCQImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffers.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type *None*

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.


```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None,
                  with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True,
                  callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.

- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.

- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.

- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type Optional[List[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[*str*, Any]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List*[Any]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (*Union[numpy.ndarray, List[Any]]*) – observations
- **action** (*Union[numpy.ndarray, List[Any]]*) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type *Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]*

sample_action (*x*)

BCQ does not support sampling action.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) –

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.BEAR

```
class d3rlpy.algos.BEAR(*, actor_learning_rate=0.0001, critic_learning_rate=0.0003, im-
    itator_learning_rate=0.0003, temp_learning_rate=0.0001, al-
    pha_learning_rate=0.001, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls=
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), imita-
    tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), al-
    pha_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), ac-
    tor_encoder_factory='default', critic_encoder_factory='default', imita-
    tor_encoder_factory='default', q_func_factory='mean', batch_size=256,
    n_frames=1, n_steps=1, gamma=0.99, tau=0.005, n_critics=2, ini-
    tial_temperature=1.0, initial_alpha=1.0, alpha_threshold=0.05,
    lam=0.75, n_action_samples=100, n_target_samples=10,
    n_mmd_action_samples=4, mmd_kernel='laplacian', mmd_sigma=20.0,
    vae_kl_weight=0.5, warmup_steps=40000, use_gpu=False, scaler=None,
    action_scaler=None, impl=None, **kwargs)
```

Bootstrapping Error Accumulation Reduction algorithm.

BEAR is a SAC-based data-driven deep reinforcement learning algorithm.

BEAR constrains the support of the policy function within data distribution by minimizing Maximum Mean Discrepancy (MMD) between the policy function and the approximated behavior policy function $\pi_\beta(a|s)$ which is optimized through L2 loss.

$$L(\beta) = \mathbb{E}_{s_t, a_t \sim D, a \sim \pi_\beta(\cdot|s_t)} [(a - a_t)^2]$$

The policy objective is a combination of SAC's objective and MMD penalty.

$$J(\phi) = J_{SAC}(\phi) - \mathbb{E}_{s_t \sim D} \alpha (\text{MMD}(\pi_\beta(\cdot|s_t), \pi_\phi(\cdot|s_t)) - \epsilon)$$

where MMD is computed as follows.

$$\text{MMD}(x, y) = \frac{1}{N^2} \sum_{i, i'} k(x_i, x_{i'}) - \frac{2}{NM} \sum_{i, j} k(x_i, y_j) + \frac{1}{M^2} \sum_{j, j'} k(y_j, y_{j'})$$

where $k(x, y)$ is a gaussian kernel $k(x, y) = \exp(-(x - y)^2 / (2\sigma^2))$.

α is also adjustable through dual gradient descent where α becomes smaller if MMD is smaller than the threshold ϵ .

References

- Kumar et al., Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **imitator_learning_rate** (*float*) – learning rate for behavior policy function.

- **temp_learning_rate** (*float*) – learning rate for temperature parameter.
- **alpha_learning_rate** (*float*) – learning rate for α .
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.
- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **imitator_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the behavior policy.
- **temp_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the temperature.
- **alpha_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for α .
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the critic.
- **imitator_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the behavior policy.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **initial_temperature** (*float*) – initial temperature value.
- **initial_alpha** (*float*) – initial α value.
- **alpha_threshold** (*float*) – threshold value described as ϵ .
- **lam** (*float*) – weight for critic ensemble.
- **n_action_samples** (*int*) – the number of action samples to compute the best action.
- **n_target_samples** (*int*) – the number of action samples to compute BCQ-like target value.
- **n_mmd_action_samples** (*int*) – the number of action samples to compute MMD.
- **mmd_kernel** (*str*) – MMD kernel function. The available options are ['gaussian', 'laplacian'].
- **mmd_sigma** (*float*) – σ for gaussian kernel in MMD calculation.
- **vae_kl_weight** (*float*) – constant weight to scale KL term for behavior policy training.
- **warmup_steps** (*int*) – the number of steps to warmup the policy function.

- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device iD or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.bear_impl.BEARImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffer.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)
Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type `List[Tuple[int, Dict[str, float]]]`

fit_batch_online (*env*, *buffer*=None, *explorer*=None, *n_epochs*=1000, *n_steps_per_epoch*=1000, *n_updates_per_epoch*=1000, *eval_interval*=10, *eval_env*=None, *eval_epsilon*=0.0, *save_metrics*=True, *save_interval*=1, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *timelimit_aware*=True, *callback*=None)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.

- **buffer** (*Optional*[*d3rlpy.online.buffer*.*Buffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional*[*Callable*[[*d3rlpy.online.iterators.AlgoProtocol*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fitter (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(episodes):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union*[*List*[*d3rlpy.dataset.Episode*], *d3rlpy.dataset.MDPDataset*]) – list of episodes to train.
- **n_epochs** (*Optional*[*int*]) – the number of epochs to train.

- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [`d3rlpy.dataset.Transition`]) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type `Dict[str, Any]`

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
```

(continues on next page)

(continued from previous page)

```
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters \mathbf{x} (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (x , $action$, $with_std=False$)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- \mathbf{x} (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (x)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters \mathbf{x} (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model ($fname$)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type `None`

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (*d3rlpy.dataset.TransitionMiniBatch*) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.CRR

```
class d3rlpy.algos.CRR(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam', betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam', betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), actor_encoder_factory='default', critic_encoder_factory='default', q_func_factory='mean', batch_size=100, n_frames=1, n_steps=1, gamma=0.99, beta=1.0, n_action_samples=4, advantage_type='mean', weight_type='exp', max_weight=20.0, n_critics=1, target_update_interval=100, target_reduction_type='min', update_actor_interval=1, use_gpu=False, scaler=None, action_scaler=None, impl=None, **kwargs)
```

Critic Regularized Regression algorithm.

CRR is a simple offline RL method similar to AWAC.

The policy is trained as a supervised regression.

$$J(\phi) = \mathbb{E}_{s_t, a_t \sim D} [\log \pi_\phi(a_t | s_t) f(Q_\theta, \pi_\phi, s_t, a_t)]$$

where f is a filter function which has several options. The first option is `binary` function.

$$f := \mathbb{I}[A_\theta(s, a) > 0]$$

The other is `exp` function.

$$f := \exp(A(s, a) / \beta)$$

The $A(s, a)$ is an average function which also has several options. The first option is `mean`.

$$A(s, a) = Q_\theta(s, a) - \frac{1}{m} \sum_j^m Q(s, a_j)$$

The other one is `max`.

$$A(s, a) = Q_\theta(s, a) - \max_j^m Q(s, a_j)$$

where $a_j \sim \pi_\phi(s)$.

In evaluation, the action is determined by Critic Weighted Policy (CWP). In CWP, the several actions are sampled from the policy function, and the final action is re-sampled from the estimated action-value distribution.

References

- Wang et al., Critic Regularized Regression.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.

- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the critic.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **beta** (*float*) – temperature value defined as β above.
- **n_action_samples** (*int*) – the number of sampled actions to calculate $A(s, a)$ and for CWP.
- **advantage_type** (*str*) – advantage function type. The available options are ['mean', 'max'].
- **weight_type** (*str*) – filter function type. The available options are ['binary', 'exp'].
- **max_weight** (*float*) – maximum weight for cross-entropy loss.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (*int*) – interval to update policy function.
- **use_gpu** (*bool, int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.crr_impl.CRRImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.Buffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_steps** (`int`) – the number of total steps to train.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **timelimit_aware** (`bool`) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (`Sequence[int]`) – observation shape.
- **action_size** (`int`) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.

- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.

- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.

- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type *Dict*[*str*, *Any*]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List*[*Any*]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action(x)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model(fname)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (`str`) – destination file path.

Return type `None`

save_params(logger)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy(fname, as_onnx=False)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.CQL

```
class d3rlpy.algos.CQL(*, actor_learning_rate=0.0001, critic_learning_rate=0.0003,
    temp_learning_rate=0.0001, alpha_learning_rate=0.0001, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), alpha_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), actor_encoder_factory='default', critic_encoder_factory='default',
    q_func_factory='mean', batch_size=256, n_frames=1, n_steps=1, gamma=0.99, tau=0.005, n_critics=2, target_reduction_type='min',
    initial_temperature=1.0, initial_alpha=1.0, alpha_threshold=10.0, conservative_weight=5.0, n_action_samples=10, soft_q_backup=False,
    use_gpu=False, scaler=None, action_scaler=None, impl=None, **kwargs)
```

Conservative Q-Learning algorithm.

CQL is a SAC-based data-driven deep reinforcement learning algorithm, which achieves state-of-the-art performance in offline RL problems.

CQL mitigates overestimation error by minimizing action-values under the current policy and maximizing values under data distribution for underestimation issue.

$$L(\theta_i) = \alpha \mathbb{E}_{s_t \sim D} \left[\log \sum_a \exp Q_{\theta_i}(s_t, a) - \mathbb{E}_{a \sim D} [Q_{\theta_i}(s, a)] - \tau \right] + L_{\text{SAC}}(\theta_i)$$

where α is an automatically adjustable value via Lagrangian dual gradient descent and τ is a threshold value. If the action-value difference is smaller than τ , the α will become smaller. Otherwise, the α will become larger to aggressively penalize action-values.

In continuous control, $\log \sum_a \exp Q(s, a)$ is computed as follows.

$$\log \sum_a \exp Q(s, a) \approx \log \left(\frac{1}{2N} \sum_{a_i \sim \text{Unif}(a)} \left[\frac{\exp Q(s, a_i)}{\text{Unif}(a)} \right] + \frac{1}{2N} \sum_{a_i \sim \pi_\phi(a|s)} \left[\frac{\exp Q(s, a_i)}{\pi_\phi(a_i|s)} \right] \right)$$

where N is the number of sampled actions.

The rest of optimization is exactly same as `d3rlpy.algos.SAC`.

References

- Kumar et al., Conservative Q-Learning for Offline Reinforcement Learning.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **temp_learning_rate** (*float*) – learning rate for temperature parameter of SAC.
- **alpha_learning_rate** (*float*) – learning rate for α .
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.

- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **temp_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the temperature.
- **alpha_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for α .
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the critic.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **tau** (`float`) – target network synchronization coefficient.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_reduction_type** (`str`) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **initial_temperature** (`float`) – initial temperature value.
- **initial_alpha** (`float`) – initial α value.
- **alpha_threshold** (`float`) – threshold value described as τ .
- **conservative_weight** (`float`) – constant weight to scale conservative loss.
- **n_action_samples** (`int`) – the number of sampled actions to compute $\log \sum_a \exp Q(s, a)$.
- **soft_q_backup** (`bool`) – flag to use SAC-style backup.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or `str`) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or `str`) – action preprocessor. The available options are ['min_max'].
- **impl** (`d3rlpy.algos.torch.cql_impl.CQLImpl`) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.Buffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_steps** (`int`) – the number of total steps to train.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **timelimit_aware** (`bool`) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (`Sequence[int]`) – observation shape.
- **action_size** (`int`) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[[List\[d3rlpy.dataset.Episode\]](#), [d3rlpy.dataset.MDPDataset](#)]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[[List\[d3rlpy.dataset.Episode\]](#)]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[[Dict\[str, Callable\[\[Any, \[List\\[d3rlpy.dataset.Episode\\]\]\(#\)\], float\]\]\]](#)*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[[d3rlpy.base.LearnableBase](#), [int](#), [int](#)], [None](#)]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type `List[Tuple[int, Dict[str, float]]]`

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-  
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-  
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*[d3rlpy.envs.batch.BatchEnv](#)*) – gym-like environment.
- **buffer** (*Optional[[d3rlpy.online.buffers.BatchBuffer](#)]*) – replay buffer.
- **explorer** (*Optional[[d3rlpy.online.explorers.Explorer](#)]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.

- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional*[*Callable*[[*d3rlpy.online.iterators.AlgoProtocol*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffer.Buffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.

- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional*[*Callable*[[*d3rlpy.online.iterators.AlgoProtocol*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(episodes):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union*[*List*[*d3rlpy.dataset.Episode*], *d3rlpy.dataset.MDPDataset*]) – list of episodes to train.
- **n_epochs** (*Optional*[*int*]) – the number of epochs to train.
- **n_steps** (*Optional*[*int*]) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.

- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List* [*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type *Dict*[*str*, *Any*]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List* [*Any*]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.AWR

```
class d3rlpy.algos.AWR(*, actor_learning_rate=5e-05, critic_learning_rate=0.0001, actor_optim_factory=d3rlpy.models.optimizers.SGDFactory(optim_cls='SGD', momentum=0.9, dampening=0, weight_decay=0, nesterov=False), critic_optim_factory=d3rlpy.models.optimizers.SGDFactory(optim_cls='SGD', momentum=0.9, dampening=0, weight_decay=0, nesterov=False), actor_encoder_factory='default', critic_encoder_factory='default', batch_size=2048, n_frames=1, gamma=0.99, batch_size_per_update=256, n_actor_updates=1000, n_critic_updates=200, lam=0.95, beta=1.0, max_weight=20.0, use_gpu=False, scaler=None, action_scaler=None, impl=None, **kwargs)
```

Advantage-Weighted Regression algorithm.

AWR is an actor-critic algorithm that trains via supervised regression way, and has shown strong performance in online and offline settings.

The value function is trained as a supervised regression problem.

$$L(\theta) = \mathbb{E}_{s_t, R_t \sim D} [(R_t - V(s_t|\theta))^2]$$

where R_t is approximated using $TD(\lambda)$ to mitigate high variance issue.

The policy function is also trained as a supervised regression problem.

$$J(\phi) = \mathbb{E}_{s_t, a_t, R_t \sim D} [\log \pi(a_t|s_t, \phi) \exp(\frac{1}{B}(R_t - V(s_t|\theta)))]$$

where B is a constant factor.

References

- [Peng et al., Advantage-Weighted Regression: Simple and Scalable Off-Policy Reinforcement Learning](#)

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for value function.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the critic.
- **batch_size** (*int*) – batch size per iteration.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **gamma** (*float*) – discount factor.
- **batch_size_per_update** (*int*) – mini-batch size.
- **n_actor_updates** (*int*) – actor gradient steps per iteration.

- **n_critic_updates** (*int*) – critic gradient steps per iteration.
- **lam** (*float*) – λ for TD(λ).
- **beta** (*float*) – B for weight scale.
- **max_weight** (*float*) – w_{\max} for weight clipping.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.awr_impl.AWRImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffers.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence*[*int*]) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type *None*

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union*[*List*[*d3rlpy.dataset.Episode*], *d3rlpy.dataset.MDPDataset*]) – list of episodes to train.
- **n_epochs** (*Optional*[*int*]) – the number of epochs to train.
- **n_steps** (*Optional*[*int*]) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {*class name*}_{*timestamp*}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional*[*List*[*d3rlpy.dataset.Episode*]]) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional*[*Dict*[*str*, *Callable*[[*Any*, *List*[*d3rlpy.dataset.Episode*]], *float*]]]) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type *List*[*Tuple*[*int*, *Dict*[*str*, *float*]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None,
                  with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True,
                  callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.


```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, **args*, ***kwargs*)

Returns predicted state values.

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations.
- **args** (*Any*) –
- **kwargs** (*Any*) –

Returns predicted state values.

Return type `numpy.ndarray`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.AWAC

```
class d3rlpy.algos.AWAC(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003, ac-
    tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0.0001, amsgrad=False),
    critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    actor_encoder_factory='default', critic_encoder_factory='default',
    q_func_factory='mean', batch_size=1024, n_frames=1, n_steps=1,
    gamma=0.99, tau=0.005, lam=1.0, n_action_samples=1,
    max_weight=20.0, n_critics=2, target_reduction_type='min',
    update_actor_interval=1, use_gpu=False, scaler=None, ac-
    tion_scaler=None, impl=None, **kwargs)
```

Advantage Weighted Actor-Critic algorithm.

AWAC is a TD3-based actor-critic algorithm that enables efficient fine-tuning where the policy is trained with offline datasets and is deployed to online training.

The policy is trained as a supervised regression.

$$J(\phi) = \mathbb{E}_{s_t, a_t \sim D} [\log \pi_{\phi}(a_t | s_t) \exp(\frac{1}{\lambda} A^{\pi}(s_t, a_t))]$$

where $A^{\pi}(s_t, a_t) = Q_{\theta}(s_t, a_t) - Q_{\theta}(s_t, a'_t)$ and $a'_t \sim \pi_{\phi}(\cdot | s_t)$

The key difference from AWR is that AWAC uses Q-function trained via TD learning for the better sample-efficiency.

References

- Nair et al., Accelerating Online Reinforcement Learning with Offline Datasets.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the critic.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory or str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.

- **lam** (*float*) – λ for weight calculation.
- **n_action_samples** (*int*) – the number of sampled actions to calculate $A^\pi(s_t, a_t)$.
- **max_weight** (*float*) – maximum weight for cross-entropy loss.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (*int*) – interval to update policy function.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.awac_impl.AWACImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If buffer is not given, ReplayBuffer will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)
Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.

- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[*int*, Dict[*str*, *float*]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffers.BatchBuffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

- **callback** (*Optional*[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]) – callable function that takes (algo, epoch, total_step) , which is called at the end of epochs.

Return type None

fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000, update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional*[d3rlpy.online.buffer.Buffer]) – replay buffer.
- **explorer** (*Optional*[d3rlpy.online.explorers.Explorer]) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional*[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]) – callable function that takes (algo, epoch, total_step) , which is called at the end of epochs.

Return type None


```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
       experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
       show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

```
classmethod from_json(fname, use_gpu=False)
```

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type *Optional[List[d3rlpy.dataset.Transition]]*

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type d3rlpy.constants.ActionSpace

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters *fname* (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters *x* (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- *x* (*Union[numpy.ndarray, List[Any]]*) – observations
- *action* (*Union[numpy.ndarray, List[Any]]*) – actions
- *with_std* (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type *Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]*

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.**Parameters** *x* (*Union[numpy.ndarray, List[Any]]*) – observations.**Returns** sampled actions.**Return type** *numpy.ndarray***save_model** (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters *fname* (*str*) – destination file path.**Return type** *None***save_params** (*logger*)

Saves configurations as params.json.

Parameters *logger* (*d3rlpy.logger.D3RLPyLogger*) – logger object.**Return type** *None***save_policy** (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- *fname* (*str*) – destination file path.
- *as_onnx* (*bool*) – flag to save as ONNX format.

Return type *None***set_params** (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters `params` (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.PLAS

```
class d3rlpy.algos.PLAS(*, actor_learning_rate=0.0001, critic_learning_rate=0.001, imita-
tor_learning_rate=0.0001, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls=
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), imita-
tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False), ac-
tor_encoder_factory='default', critic_encoder_factory='default', imita-
tor_encoder_factory='default', q_func_factory='mean', batch_size=100,
n_frames=1, n_steps=1, gamma=0.99, tau=0.005, n_critics=2,
target_reduction_type='mix', update_actor_interval=1, lam=0.75,
warmup_steps=500000, beta=0.5, use_gpu=False, scaler=None, ac-
tion_scaler=None, impl=None, **kwargs)
```

Policy in Latent Action Space algorithm.

PLAS is an offline deep reinforcement learning algorithm whose policy function is trained in latent space of Conditional VAE. Unlike other algorithms, PLAS can achieve good performance by using its less constrained policy function.

$$a \sim p_{\beta}(a|s, z = \pi_{\phi}(s))$$

where β is a parameter of the decoder in Conditional VAE.

References

- Zhou et al., PLAS: latent action space for offline reinforcement learning.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **imitator_learning_rate** (*float*) – learning rate for Conditional VAE.
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.

- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **imitator_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the conditional VAE.
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the critic.
- **imitator_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the conditional VAE.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **tau** (`float`) – target network synchronization coefficient.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_reduction_type** (`str`) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (`int`) – interval to update policy function.
- **lam** (`float`) – weight factor for critic ensemble.
- **warmup_steps** (`int`) – the number of steps to warmup the VAE.
- **beta** (`float`) – KL regularization term for Conditional VAE.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or `str`) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or `str`) – action preprocessor. The available options are ['min_max'].
- **impl** (`d3rlpy.algos.torch.bcq_impl.BCQImpl`) – algorithm implementation.

Methods

build_with_dataset (`dataset`)

Instantiate implementation object with MDPDataset object.

Parameters `dataset` (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (`env`)

Instantiate implementation object with OpenAI Gym object.

Parameters `env` (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffer.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.

- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None,
                  with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.

- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Return type *None*

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.

- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type Optional[List[d3rlpy.dataset.Transition]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type d3rlpy.constants.ActionSpace

get_params (*deep=True*)
Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters *deep* (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)
Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters *fname* (*str*) – source file path.

Return type None

predict (*x*)
Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters *x* (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (*x, action, with_std=False*)
Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
```

(continues on next page)

(continued from previous page)

```
# values.shape == (100,)
values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values**Return type** `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`**sample_action** (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.**Parameters** **x** (`Union[numpy.ndarray, List[Any]]`) – observations.**Returns** sampled actions.**Return type** `numpy.ndarray`**save_model** (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.**Return type** `None`**save_params** (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.**Return type** `None`**save_policy** (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).

- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.PLASWithPerturbation

```
class d3rlpy.algos.PLASWithPerturbation(*,
                                       actor_learning_rate=0.0001,
                                       critic_learning_rate=0.001,
                                       imitator_learning_rate=0.0001,
                                       actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Ada
                                       betas=(0.9, 0.999), eps=1e-08,
                                       weight_decay=0, amsgrad=False),
                                       critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='A
                                       betas=(0.9, 0.999), eps=1e-08,
                                       weight_decay=0, amsgrad=False), imita
                                       tor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Ada
                                       betas=(0.9, 0.999), eps=1e-08, weight_decay=0,
                                       amsgrad=False), actor_encoder_factory='default',
                                       critic_encoder_factory='default', imi
                                       tator_encoder_factory='default',
                                       q_func_factory='mean', batch_size=100,
                                       n_frames=1, n_steps=1, gamma=0.99, tau=0.005,
                                       n_critics=2, target_reduction_type='mix',
                                       update_actor_interval=1, lam=0.75, ac
                                       tion_flexibility=0.05, warmup_steps=500000,
                                       beta=0.5, use_gpu=False, scaler=None, ac
                                       tion_scaler=None, impl=None, **kwargs)
```

Policy in Latent Action Space algorithm with perturbation layer.

PLAS with perturbation layer enables PLAS to output out-of-distribution action.

References

- Zhou et al., PLAS: latent action space for offline reinforcement learning.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **imitator_learning_rate** (*float*) – learning rate for Conditional VAE.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **imitator_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the conditional VAE.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the critic.
- **imitator_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the conditional VAE.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.

- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (*int*) – interval to update policy function.
- **lam** (*float*) – weight factor for critic ensemble.
- **action_flexibility** (*float*) – output scale of perturbation layer.
- **warmup_steps** (*int*) – the number of steps to warmup the VAE.
- **beta** (*float*) – KL regularization term for Conditional VAE.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.bcq_impl.BCQImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.

- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.

- **scorers** (*Optional*[*Dict*[*str*, *Callable*[[*Any*, *List*[*d3rlpy.dataset.Episode*]], *float*]]]) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type *List*[*Tuple*[*int*, *Dict*[*str*, *float*]]]

fit_batch_online (*env*, *buffer=None*, *explorer=None*, *n_epochs=1000*, *n_steps_per_epoch=1000*, *n_updates_per_epoch=1000*, *eval_interval=10*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffers.BatchBuffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.

- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fit_online (*env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000, update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If `None`, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If `False`, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type `Generator[Tuple[int, Dict[str, float]], None, None]`

classmethod from_json (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type Optional[List[d3rlpy.dataset.Transition]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type d3rlpy.constants.ActionSpace

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (fname)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (str) – source file path.

Return type None

predict (x)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (Union[numpy.ndarray, List[Any]]) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (x, action, with_std=False)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (Union[numpy.ndarray, List[Any]]) – observations
- **action** (Union[numpy.ndarray, List[Any]]) – actions
- **with_std** (bool) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters *x* (Union[numpy.ndarray, List[Any]]) – observations.

Returns sampled actions.

Return type numpy.ndarray

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters *fname* (*str*) – destination file path.

Return type None

save_params (*logger*)

Saves configurations as params.json.

Parameters *logger* (d3rlpy.logger.D3RLPyLogger) – logger object.

Return type None

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- *fname* (*str*) – destination file path.
- *as_onnx* (*bool*) – flag to save as ONNX format.

Return type None

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type d3rlpy.base.LearnableBase

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (d3rlpy.dataset.TransitionMiniBatch) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.MOPO

```
class d3rlpy.algos.MOPO(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003,
    temp_learning_rate=0.0003, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    actor_encoder_factory='default', critic_encoder_factory='default',
    q_func_factory='mean', batch_size=100, n_frames=1, n_steps=1,
    gamma=0.99, tau=0.005, n_critics=2, target_reduction_type='min',
    update_actor_interval=1, initial_temperature=1.0, dynamics=None,
    rollout_interval=1000, rollout_horizon=5, rollout_batch_size=50000,
    lam=1.0, real_ratio=0.05, generated_maxlen=1250000, use_gpu=False,
    scaler=None, action_scaler=None, impl=None, **kwargs)
```

Model-based Offline Policy Optimization.

MOPO is a model-based RL approach for offline policy optimization. MOPO leverages the probabilistic ensemble dynamics model to generate new dynamics data with uncertainty penalties. The ensemble dynamics model consists of N probabilistic models $\{T_{\theta_i}\}_{i=1}^N$. At each epoch, new transitions are generated via randomly picked dynamics model T_{θ} .

$$s_{t+1}, r_{t+1} \sim T_{\theta}(s_t, a_t)$$

where $s_t \sim D$ for the first step, otherwise s_t is the previous generated observation, and $a_t \sim \pi(\cdot|s_t)$. The generated r_{t+1} would be far from the ground truth if the actions sampled from the policy function is out-of-distribution. Thus, the uncertainty penalty regularizes this bias.

$$\tilde{r}_{t+1} = r_{t+1} - \lambda \max_{i=1}^N \|\Sigma_i(s_t, a_t)\|$$

where $\Sigma(s_t, a_t)$ is the estimated variance. Finally, the generated transitions $(s_t, a_t, r_{t+1}, s_{t+1})$ are appended to dataset D . This generation process starts with randomly sampled `n_initial_transitions` transitions till horizon steps.

Note: Currently, MOPO only supports vector observations.

References

- Yu et al., MOPO: Model-based Offline Policy Optimization.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **temp_learning_rate** (*float*) – learning rate for temperature parameter.
- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.
- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **temp_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the temperature.
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory for the critic.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **tau** (*float*) – target network synchronization coefficient.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (*int*) – interval to update policy function.
- **initial_temperature** (*float*) – initial temperature value.
- **dynamics** (`d3rlpy.dynamics.DynamicsBase`) – dynamics object.
- **rollout_interval** (*int*) – the number of steps before rollout.
- **rollout_horizon** (*int*) – the rollout step length.
- **rollout_batch_size** (*int*) – the number of initial transitions for rollout.

- **lam** (*float*) – λ for uncertainty penalties.
- **real_ratio** (*float*) – the real of dataset samples in a mini-batch.
- **generated_maxlen** (*int*) – the maximum number of generated samples.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessing.Scaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are ['min_max'].
- **impl** (*d3rlpy.algos.torch.sac_impl.SACImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffers.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)
Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type `List[Tuple[int, Dict[str, float]]]`

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  

n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  

eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  

show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.


```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type *None*

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.COMBO

```
class d3rlpy.algos.COMBO(*,      actor_learning_rate=0.0001,      critic_learning_rate=0.0003,
                           temp_learning_rate=0.0001, actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_
                           betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
                           critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                           betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
                           temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                           betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
                           actor_encoder_factory='default',      critic_encoder_factory='default',
                           q_func_factory='mean', batch_size=256, n_frames=1, n_steps=1,
                           gamma=0.99, tau=0.005, n_critics=2, target_reduction_type='min',
                           update_actor_interval=1,      initial_temperature=1.0,      conserva-
                           tive_weight=1.0,      n_action_samples=10,      soft_q_backup=False,
                           dynamics=None, rollout_interval=1000, rollout_horizon=5, roll-
                           out_batch_size=50000, real_ratio=0.5, generated_maxlen=1250000,
                           use_gpu=False, scaler=None, action_scaler=None, impl=None,
                           **kwargs)
```

Conservative Offline Model-Based Optimization.

COMBO is a model-based RL approach for offline policy optimization. COMBO is similar to MOPO, but it also leverages conservative loss proposed in CQL.

$$L(\theta_i) = \mathbb{E}_{s \sim d_M} [\log \sum_a \exp Q_{\theta_i}(s, a)] - \mathbb{E}_{s, a \sim D} [Q_{\theta_i}(s, a)] + L_{\text{SAC}}(\theta_i)$$

Note: Currently, COMBO only supports vector observations.

References

- Yu et al., COMBO: Conservative Offline Model-Based Policy Optimization.

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **temp_learning_rate** (*float*) – learning rate for temperature parameter.

- **actor_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the actor.
- **critic_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the critic.
- **temp_optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory for the temperature.
- **actor_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the actor.
- **critic_encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory for the critic.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **tau** (`float`) – target network synchronization coefficient.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_reduction_type** (`str`) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **update_actor_interval** (`int`) – interval to update policy function.
- **initial_temperature** (`float`) – initial temperature value.
- **conservative_weight** (`float`) – constant weight to scale conservative loss.
- **n_action_samples** (`int`) – the number of sampled actions to compute $\log \sum_a \exp Q(s, a)$.
- **soft_q_backup** (`bool`) – flag to use SAC-style backup.
- **dynamics** (`d3rlpy.dynamics.DynamicsBase`) – dynamics object.
- **rollout_interval** (`int`) – the number of steps before rollout.
- **rollout_horizon** (`int`) – the rollout step length.
- **rollout_batch_size** (`int`) – the number of initial transitions for rollout.
- **real_ratio** (`float`) – the real of dataset samples in a mini-batch.
- **generated_maxlen** (`int`) – the maximum number of generated samples.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessing.Scaler` or `str`) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or `str`) – action preprocessor. The available options are ['min_max'].
- **impl** (`d3rlpy.algos.torch.combo_impl.COMBOImpl`) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.Buffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_steps** (`int`) – the number of total steps to train.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **timelimit_aware** (`bool`) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (`Sequence[int]`) – observation shape.
- **action_size** (`int`) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.

- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional*[*Callable*[[*d3rlpy.online.iterators.AlgoProtocol*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffer.Buffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.

- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.

- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type *Dict*[*str*, *Any*]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List*[*Any*]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.RandomPolicy

```
class d3rlpy.algos.RandomPolicy (*,      distribution='uniform',      normal_std=1.0,      ac-
                                tion_scaler=None, **kwargs)
```

Random Policy for continuous control algorithm.

This is designed for data collection and lightweight interaction tests. `fit` and `fit_online` methods will raise exceptions.

Parameters

- **distribution** (*str*) – random distribution. The available options are `['uniform', 'normal']`.
- **normal_std** (*float*) – standard deviation of the normal distribution. This is only used when `distribution='normal'`.
- **action_scaler** (*d3rlpy.preprocessing.ActionScaler* or *str*) – action preprocessor. The available options are `['min_max']`.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with `MDPDataSet` object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataSet*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If `buffer` is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffers.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None,  
                  with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type None

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type *None*

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

3.1.2 Discrete control algorithms

<code>d3rlpy.algos.DiscreteBC</code>	Behavior Cloning algorithm for discrete control.
<code>d3rlpy.algos.DQN</code>	Deep Q-Network algorithm.
<code>d3rlpy.algos.DoubleDQN</code>	Double Deep Q-Network algorithm.
<code>d3rlpy.algos.DiscreteSAC</code>	Soft Actor-Critic algorithm for discrete action-space.
<code>d3rlpy.algos.DiscreteBCQ</code>	Discrete version of Batch-Constrained Q-learning algorithm.
<code>d3rlpy.algos.DiscreteCQL</code>	Discrete version of Conservative Q-Learning algorithm.
<code>d3rlpy.algos.DiscreteAWR</code>	Discrete version of Advantage-Weighted Regression algorithm.
<code>d3rlpy.algos.DiscreteRandomPolicy</code>	Random Policy for discrete control algorithm.

`d3rlpy.algos.DiscreteBC`

```
class d3rlpy.algos.DiscreteBC(* , learning_rate=0.001, optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls=  
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, ams-  
    grad=False), encoder_factory='default', batch_size=100,  
    n_frames=1, beta=0.5, use_gpu=False, scaler=None,  
    impl=None, **kwargs)
```

Behavior Cloning algorithm for discrete control.

Behavior Cloning (BC) is to imitate actions in the dataset via a supervised learning approach. Since BC is only imitating action distributions, the performance will be close to the mean of the dataset even though BC mostly works better than online RL algorithms.

$$L(\theta) = \mathbb{E}_{a_t, s_t \sim D} \left[- \sum_a p(a|s_t) \log \pi_{\theta}(a|s_t) \right]$$

where $p(a|s_t)$ is implemented as a one-hot vector.

Parameters

- **learning_rate** (`float`) – learning rate.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory.

- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **beta** (*float*) – regularization factor.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (*d3rlpy.algos.torch.bc_impl.DiscreteBCImpl*) – implementation of the algorithm.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type *d3rlpy.online.buffer.Buffer*

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.

- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type `List[Tuple[int, Dict[str, float]]]`

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  

n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  

eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  

show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.


```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

value prediction is not supported by BC algorithms.

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) –
- **action** (`Union[numpy.ndarray, List[Any]]`) –
- **with_std** (*bool*) –

Return type `numpy.ndarray`

sample_action (*x*)

sampling action is not supported by BC algorithm.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) –

Return type `None`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DQN

```
class d3rlpy.algos.DQN(*, learning_rate=6.25e-05, optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
    encoder_factory='default', q_func_factory='mean', batch_size=32,
    n_frames=1, n_steps=1, gamma=0.99, n_critics=1, target_reduction_type='min', target_update_interval=8000, use_gpu=False,
    scaler=None, impl=None, **kwargs)
```

Deep Q-Network algorithm.

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} [(r_{t+1} + \gamma \max_a Q_{\theta'}(s_{t+1}, a) - Q_{\theta}(s_t, a_t))^2]$$

where θ' is the target network parameter. The target network parameter is synchronized every `target_update_interval` iterations.

References

- Mnih et al., Human-level control through deep reinforcement learning.

Parameters

- **learning_rate** (*float*) – learning rate.
- **optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory* or *str*) – optimizer factory.
- **encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **target_update_interval** (*int*) – interval to update the target network.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (*d3rlpy.algos.torch.dqn_impl.DQNImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.Buffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_steps** (`int`) – the number of total steps to train.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **timelimit_aware** (`bool`) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (`Sequence[int]`) – observation shape.
- **action_size** (`int`) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.

- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.

- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(episodes):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.

- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type *Dict*[*str*, *Any*]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List*[*Any*]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DoubleDQN

```
class d3rlpy.algos.DoubleDQN(*,
                             learning_rate=6.25e-05,
                             optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                         betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
                             encoder_factory='default', q_func_factory='mean', batch_size=32,
                             n_frames=1, n_steps=1, gamma=0.99, n_critics=1, target_reduction_type='min',
                             target_update_interval=8000,
                             use_gpu=False, scaler=None, impl=None, **kwargs)
```

Double Deep Q-Network algorithm.

The difference from DQN is that the action is taken from the current Q function instead of the target Q function. This modification significantly decreases overestimation bias of TD learning.

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} [(r_{t+1} + \gamma Q_{\theta'}(s_{t+1}, \operatorname{argmax}_a Q_{\theta}(s_{t+1}, a)) - Q_{\theta}(s_t, a_t))^2]$$

where θ' is the target network parameter. The target network parameter is synchronized every *target_update_interval* iterations.

References

- Hasselt et al., Deep reinforcement learning with double Q-learning.

Parameters

- **learning_rate** (*float*) – learning rate.
- **optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory.
- **encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **n_critics** (*int*) – the number of Q functions.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **target_update_interval** (*int*) – interval to synchronize the target network.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (*d3rlpy.algos.torch.dqn_impl.DoubleDQNImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.Buffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_steps** (`int`) – the number of total steps to train.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **timelimit_aware** (`bool`) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (`Sequence[int]`) – observation shape.
- **action_size** (`int`) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(epochs, n_steps=1000000)
```

Parameters

- **dataset** (*Union[[List\[d3rlpy.dataset.Episode\]](#), [d3rlpy.dataset.MDPDataset](#)]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[[List\[d3rlpy.dataset.Episode\]](#)]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, [List\[d3rlpy.dataset.Episode\]](#)], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[[d3rlpy.base.LearnableBase](#), int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type `List[Tuple[int, Dict[str, float]]]`

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*[d3rlpy.envs.batch.BatchEnv](#)*) – gym-like environment.
- **buffer** (*Optional[[d3rlpy.online.buffers.BatchBuffer](#)]*) – replay buffer.
- **explorer** (*Optional[[d3rlpy.online.explorers.Explorer](#)]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.

- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.

- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.

- **total_step** (*int*) – the total update steps.
- **transitions** (*List* [*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List* [*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type *Dict*[*str*, *Any*]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (*Union*[*numpy.ndarray*, *List* [*Any*]]) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch, total_step, batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DiscreteSAC

```
class d3rlpy.algos.DiscreteSAC(*, actor_learning_rate=0.0003, critic_learning_rate=0.0003,
                                temp_learning_rate=0.0003,
                                actor_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                betas=(0.9, 0.999), eps=0.0001,
                                weight_decay=0, amsgrad=False),
                                critic_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                betas=(0.9, 0.999), eps=0.0001,
                                weight_decay=0, amsgrad=False),
                                temp_optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                betas=(0.9, 0.999), eps=0.0001, weight_decay=0,
                                amsgrad=False), actor_encoder_factory='default',
                                critic_encoder_factory='default', q_func_factory='mean',
                                batch_size=64, n_frames=1, n_steps=1, gamma=0.99,
                                n_critics=2, initial_temperature=1.0, tar-
                                get_update_interval=8000, use_gpu=False, scaler=None,
                                impl=None, **kwargs)
```

Soft Actor-Critic algorithm for discrete action-space.

This discrete version of SAC is built based on continuous version of SAC with additional modifications.

The target state-value is calculated as expectation of all action-values.

$$V(s_t) = \pi_\phi(s_t)^T [Q_\theta(s_t) - \alpha \log(\pi_\phi(s_t))]$$

Similarly, the objective function for the temperature parameter is as follows.

$$J(\alpha) = \pi_\phi(s_t)^T [-\alpha(\log(\pi_\phi(s_t)) + H)]$$

Finally, the objective function for the policy function is as follows.

$$J(\phi) = \mathbb{E}_{s_t \sim D} [\pi_\phi(s_t)^T [\alpha \log(\pi_\phi(s_t)) - Q_\theta(s_t)]]$$

References

- [Christodoulou, Soft Actor-Critic for Discrete Action Settings.](#)

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for Q functions.
- **temp_learning_rate** (*float*) – learning rate for temperature parameter.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **temp_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the temperature.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory or str*) – encoder factory for the actor.

- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the critic.
- **q_func_factory** (*d3rlpy.models.q_functions.QFunctionFactory* or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **initial_temperature** (*float*) – initial temperature value.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (*d3rlpy.algos.torch.sac_impl.DiscreteSACImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters *dataset* (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters *env* (*gym.core.Env*) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time_limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, *ReplayBuffer* will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn *terminal flag* *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.

- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[*int*, Dict[*str*, *float*]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional*[*d3rlpy.online.buffers.BatchBuffer*]) – replay buffer.
- **explorer** (*Optional*[*d3rlpy.online.explorers.Explorer*]) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional*[*gym.core.Env*]) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step) , which is called at the end of epochs.

Return type `None`

fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000, update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step) , which is called at the end of epochs.

Return type `None`

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
       experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
       show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when **n_steps** is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

```
classmethod from_json(fname, use_gpu=False)
```

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type `Dict[str, Any]`

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters *fname* (*str*) – source file path.

Return type *None*

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters *x* (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type *numpy.ndarray*

predict_value (*x, action, with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- *x* (*Union[numpy.ndarray, List[Any]]*) – observations
- *action* (*Union[numpy.ndarray, List[Any]]*) – actions
- *with_std* (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type *Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]*

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.**Parameters** *x* (*Union[numpy.ndarray, List[Any]]*) – observations.**Returns** sampled actions.**Return type** *numpy.ndarray***save_model** (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters *fname* (*str*) – destination file path.**Return type** *None***save_params** (*logger*)

Saves configurations as params.json.

Parameters *logger* (*d3rlpy.logger.D3RLPyLogger*) – logger object.**Return type** *None***save_policy** (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type *None***set_params** (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```


Parameters `params` (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.DiscreteBCQ

```
class d3rlpy.algos.DiscreteBCQ(*, learning_rate=6.25e-05, op-
    tim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
    betas=(0.9, 0.999), eps=1e-08, weight_decay=0,
    amsgrad=False), encoder_factory='default',
    q_func_factory='mean', batch_size=32, n_frames=1, n_steps=1,
    gamma=0.99, n_critics=1, target_reduction_type='min', ac-
    tion_flexibility=0.3, beta=0.5, target_update_interval=8000,
    use_gpu=False, scaler=None, impl=None, **kwargs)
```

Discrete version of Batch-Constrained Q-learning algorithm.

Discrete version takes theories from the continuous version, but the algorithm is much simpler than that. The imitation function $G_\omega(a|s)$ is trained as supervised learning just like Behavior Cloning.

$$L(\omega) = \mathbb{E}_{a_t, s_t \sim D} \left[- \sum_a p(a|s_t) \log G_\omega(a|s_t) \right]$$

With this imitation function, the greedy policy is defined as follows.

$$\pi(s_t) = \operatorname{argmax}_{a|G_\omega(a|s_t)/\max_{\tilde{a}} G_\omega(\tilde{a}|s_t) > \tau} Q_\theta(s_t, a)$$

which eliminates actions with probabilities τ times smaller than the maximum one.

Finally, the loss function is computed in Double DQN style with the above constrained policy.

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} [(r_{t+1} + \gamma Q_{\theta'}(s_{t+1}, \pi(s_{t+1})) - Q_\theta(s_t, a_t))^2]$$

References

- Fujimoto et al., Off-Policy Deep Reinforcement Learning without Exploration.
- Fujimoto et al., Benchmarking Batch Deep Reinforcement Learning Algorithms.

Parameters

- **learning_rate** (`float`) – learning rate.

- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **action_flexibility** (*float*) – probability threshold represented as τ .
- **beta** (*float*) – regularization term for imitation function.
- **target_update_interval** (*int*) – interval to update the target network.
- **use_gpu** (*bool*, *int* or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (`d3rlpy.algos.torch.bcq_impl.DiscreteBCQImpl`) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters **dataset** (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters **env** (`gym.core.Env`) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If buffer is not given, ReplayBuffer will be internally created.

Parameters

- **env** (`gym.core.Env`) – gym-like environment.
- **buffer** (*Optional* [`d3rlpy.online.buffers.Buffer`]) – replay buffer.

- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape, action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.

- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.

- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fit_online (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *n_steps_per_epoch=10000*, *update_interval=1*, *update_start_step=0*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

- **callback** (*Optional*[*Callable*[[*d3rlpy.online.iterators.AlgoProtocol*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fitter (dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union*[*List*[*d3rlpy.dataset.Episode*], *d3rlpy.dataset.MDPDataset*]) – list of episodes to train.
- **n_epochs** (*Optional*[*int*]) – the number of epochs to train.
- **n_steps** (*Optional*[*int*]) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when n_steps is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional*[*str*]) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional*[*str*]) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional*[*List*[*d3rlpy.dataset.Episode*]]) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional*[*Dict*[*str*, *Callable*[[*Any*, *List*[*d3rlpy.dataset.Episode*], *float*]]]) – list of scorer functions used with eval_episodes.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type *Generator*[*Tuple*[*int*, *Dict*[*str*, *float*]], *None*, *None*]

classmethod `from_json(fname, use_gpu=False)`

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (fname)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (str) – source file path.

Return type None

predict (x)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (Union[numpy.ndarray, List[Any]]) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (x, action, with_std=False)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (Union[numpy.ndarray, List[Any]]) – observations
- **action** (Union[numpy.ndarray, List[Any]]) – actions
- **with_std** (bool) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable bootstrap flag and increase n_critics value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters *x* (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters *fname* (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters *logger* (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- *fname* (*str*) – destination file path.
- *as_onnx* (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type d3rlpy.base.LearnableBase

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (d3rlpy.dataset.TransitionMiniBatch) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

d3rlpy.algos.DiscreteCQL

```
class d3rlpy.algos.DiscreteCQL(*,                               learning_rate=6.25e-05,                               op-
                               tim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                               betas=(0.9, 0.999), eps=1e-08, weight_decay=0,
                               amsgrad=False), encoder_factory='default',
                               q_func_factory='mean', batch_size=32, n_frames=1, n_steps=1,
                               gamma=0.99, n_critics=1, target_reduction_type='min', tar-
                               get_update_interval=8000, use_gpu=False, scaler=None,
                               impl=None, **kwargs)
```

Discrete version of Conservative Q-Learning algorithm.

Discrete version of CQL is a DoubleDQN-based data-driven deep reinforcement learning algorithm (the original paper uses DQN), which achieves state-of-the-art performance in offline RL problems.

CQL mitigates overestimation error by minimizing action-values under the current policy and maximizing values under data distribution for underestimation issue.

$$L(\theta) = \mathbb{E}_{s_t \sim D} [\log \sum_a \exp Q_\theta(s_t, a) - \mathbb{E}_{a \sim D} [Q_\theta(s, a)]] + L_{\text{DoubleDQN}}(\theta)$$

References

- Kumar et al., Conservative Q-Learning for Offline Reinforcement Learning.

Parameters

- **learning_rate** (*float*) – learning rate.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or *str*) – Q function factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – N-step TD calculation.
- **gamma** (*float*) – discount factor.
- **n_critics** (*int*) – the number of Q functions for ensemble.
- **target_reduction_type** (*str*) – ensemble reduction method at target value estimation. The available options are ['min', 'max', 'mean', 'mix', 'none'].
- **target_update_interval** (*int*) – interval to synchronize the target network.
- **use_gpu** (*bool*, *int* or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard']
- **impl** (`d3rlpy.algos.torch.cql_impl.DiscreteCQLImpl`) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters **dataset** (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters **env** (`gym.core.Env`) – gym-like environment.

Return type *None*

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If buffer is not given, ReplayBuffer will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape, action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.

- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.

- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Return type *None*

fit_online (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *n_steps_per_epoch=10000*, *update_interval=1*, *update_start_step=0*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.

- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Return type `None`

fitter (`dataset`, `n_epochs=None`, `n_steps=None`, `n_steps_per_epoch=10000`, `save_metrics=True`, `experiment_name=None`, `with_timestamp=True`, `logdir='d3rlpy_logs'`, `verbose=True`, `show_progress=True`, `tensorboard_dir=None`, `eval_episodes=None`, `save_interval=1`, `scorers=None`, `shuffle=True`, `callback=None`)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.

- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type *Generator*[*Tuple*[*int*, *Dict*[*str*, *float*]], *None*, *None*]

classmethod **from_json** (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional*[*Union*[*bool*, *int*, *d3rlpy.gpu.Device*]]) – flag to use GPU, device ID or device.

Returns algorithm.

Return type *d3rlpy.base.LearnableBase*

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List*[*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```

params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)

```

Parameters `deep` (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)

Load neural network parameters.

```

algo.load_model('model.pt')

```

Parameters `fname` (*str*) – source file path.

Return type None

predict (*x*)

Returns greedy actions.

```

# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control

```

Parameters `x` (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (*x, action, with_std=False*)

Returns predicted action-values.

```

# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)

```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DiscreteAWR

```
class d3rlpy.algos.DiscreteAWR(*, actor_learning_rate=5e-05, critic_learning_rate=0.0001, ac-
    tor_optim_factory=d3rlpy.models.optimizers.SGDFactory(optim_cls='SGD',
    momentum=0.9,                                         dampening=0,
    weight_decay=0,                                       nesterov=False),
    critic_optim_factory=d3rlpy.models.optimizers.SGDFactory(optim_cls='SGD',
    momentum=0.9,      dampening=0,      weight_decay=0,
    nesterov=False),      actor_encoder_factory='default',
    critic_encoder_factory='default',      batch_size=2048,
    n_frames=1,      gamma=0.99,      batch_size_per_update=256,
    n_actor_updates=1000,      n_critic_updates=200,      lam=0.95,
    beta=1.0,      max_weight=20.0,      use_gpu=False,      scaler=None,
    action_scaler=None, impl=None, **kwargs)
```

Discrete veriosn of Advantage-Weighted Regression algorithm.

AWR is an actor-critic algorithm that trains via supervised regression way, and has shown strong performance in online and offline settings.

The value function is trained as a supervised regression problem.

$$L(\theta) = \mathbb{E}_{s_t, R_t \sim D} [(R_t - V(s_t|\theta))^2]$$

where R_t is approximated using TD(λ) to mitigate high variance issue.

The policy function is also trained as a supervised regression problem.

$$J(\phi) = \mathbb{E}_{s_t, a_t, R_t \sim D} [\log \pi(a_t|s_t, \phi) \exp(\frac{1}{B}(R_t - V(s_t|\theta)))]$$

where B is a constant factor.

References

- Peng et al., Advantage-Weighted Regression: Simple and Scalable Off-Policy Reinforcement Learning

Parameters

- **actor_learning_rate** (*float*) – learning rate for policy function.
- **critic_learning_rate** (*float*) – learning rate for value function.
- **actor_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the actor.
- **critic_optim_factory** (*d3rlpy.models.optimizers.OptimizerFactory*) – optimizer factory for the critic.
- **actor_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the actor.
- **critic_encoder_factory** (*d3rlpy.models.encoders.EncoderFactory* or *str*) – encoder factory for the critic.
- **batch_size** (*int*) – batch size per iteration.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **gamma** (*float*) – discount factor.
- **batch_size_per_update** (*int*) – mini-batch size.
- **n_actor_updates** (*int*) – actor gradient steps per iteration.
- **n_critic_updates** (*int*) – critic gradient steps per iteration.
- **lam** (*float*) – λ for TD(λ).
- **beta** (*float*) – B for weight scale.
- **max_weight** (*float*) – w_{\max} for weight clipping.
- **use_gpu** (*bool*, *int* or *d3rlpy.gpu.Device*) – flag to use GPU, device ID or device.
- **scaler** (*d3rlpy.preprocessingScaler* or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **impl** (*d3rlpy.algos.torch.awr_impl.DiscreteAWRImpl*) – algorithm implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with MDPDataset object.

Parameters **dataset** (*d3rlpy.dataset.MDPDataset*) – dataset.

Return type *None*

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters **env** (*gym.core.Env*) – gym-like environment.

Return type `None`

collect (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *show_progress=True*, *time-limit_aware=True*)

Collects data via interaction with environment.

If *buffer* is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn `terminal flag False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffer.Buffer`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.

- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.

- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.

- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(episodes):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.

- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type d3rlpy.base.LearnableBase

generate_new_data (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type Optional[List[d3rlpy.dataset.Transition]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type d3rlpy.constants.ActionSpace

get_params (*deep=True*)
Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters *deep* (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)
Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters *fname* (*str*) – source file path.

Return type None

predict (*x*)
Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters *x* (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (*x, *args, **kwargs*)
Returns predicted state values.

Parameters

- *x* (*Union[numpy.ndarray, List[Any]]*) – observations.
- *args* (*Any*) –
- *kwargs* (*Any*) –

Returns predicted state values.

Return type numpy.ndarray

sample_action (*x*)
Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type `None`

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

d3rlpy.algos.DiscreteRandomPolicy

class d3rlpy.algos.DiscreteRandomPolicy (**kwargs)

Random Policy for discrete control algorithm.

This is designed for data collection and lightweight interaction tests. `fit` and `fit_online` methods will raise exceptions.

Methods

build_with_dataset (dataset)

Instantiate implementation object with MDPDataset object.

Parameters **dataset** (d3rlpy.dataset.MDPDataset) – dataset.

Return type None

build_with_env (env)

Instantiate implementation object with OpenAI Gym object.

Parameters **env** (gym.core.Env) – gym-like environment.

Return type None

collect (env, buffer=None, explorer=None, n_steps=1000000, show_progress=True, time_limit_aware=True)

Collects data via interaction with environment.

If `buffer` is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (gym.core.Env) – gym-like environment.
- **buffer** (Optional[d3rlpy.online.buffers.Buffer]) – replay buffer.
- **explorer** (Optional[d3rlpy.online.explorers.Explorer]) – action explorer.
- **n_steps** (int) – the number of total steps to train.
- **show_progress** (bool) – flag to show progress bar for iterations.
- **timelimit_aware** (bool) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type d3rlpy.online.buffers.Buffer

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type *None*

fit (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-  
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-  
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (`d3rlpy.envs.batch.BatchEnv`) – gym-like environment.
- **buffer** (`Optional[d3rlpy.online.buffers.BatchBuffer]`) – replay buffer.
- **explorer** (`Optional[d3rlpy.online.explorers.Explorer]`) – action explorer.
- **n_epochs** (`int`) – the number of epochs to train.
- **n_steps_per_epoch** (`int`) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (`int`) – the number of updates per epoch.
- **eval_interval** (`int`) – the number of epochs before evaluation.
- **eval_env** (`Optional[gym.core.Env]`) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (`float`) – ϵ -greedy factor during evaluation.
- **save_metrics** (`bool`) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (`int`) – the number of epochs before saving models.
- **experiment_name** (`Optional[str]`) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (`bool`) – flag to add timestamp string to the last of directory name.
- **logdir** (`str`) – root directory name to save logs.
- **verbose** (`bool`) – flag to show logged information on stdout.
- **show_progress** (`bool`) – flag to show progress bar for iterations.
- **tensorboard_dir** (`Optional[str]`) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (`bool`) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (`Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]`) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type None

```
fit_online(env, buffer=None, explorer=None, n_steps=1000000, n_steps_per_epoch=10000,
            update_interval=1, update_start_step=0, eval_env=None, eval_epsilon=0.0,
            save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True,
            logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None,
            timelimit_aware=True, callback=None)
```

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag False when TimeLimit.truncated flag is True, which is designed to incorporate with gym.wrappers.TimeLimit.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

```
fitter(dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True,
        experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,
        show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None,
        shuffle=True, callback=None)
```

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```

for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)

```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod from_json (*fname, use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```

from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

```

(continues on next page)

(continued from previous page)

```
# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.**Return type** `d3rlpy.base.LearnableBase`**generate_new_data** (*epoch, total_step, transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.**Return type** `Optional[List[d3rlpy.dataset.Transition]]`**get_action_type** ()

Returns action type (continuous or discrete).

Returns action type.**Return type** `d3rlpy.constants.ActionSpace`**get_params** (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.**Returns** attribute values in dictionary.**Return type** `Dict[str, Any]`**load_model** (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type `None`

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations

Returns greedy actions

Return type `numpy.ndarray`

predict_value (*x*, *action*, *with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (*bool*) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (*Union[numpy.ndarray, List[Any]]*) – observations.

Returns sampled actions.

Return type *numpy.ndarray*

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type *None*

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type *None*

save_policy (*fname, as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type *None*

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type `Optional[ImplBase]`

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

3.2 Q Functions

d3rlpy provides various Q functions including state-of-the-arts, which are internally used in algorithm objects. You can switch Q functions by passing `q_func_factory` argument at algorithm initialization.

```
from d3rlpy.algos import CQL

cql = CQL(q_func_factory='qr') # use Quantile Regression Q function
```

Also you can change hyper parameters.

```
from d3rlpy.models.q_functions import QRQFunctionFactory

q_func = QRQFunctionFactory(n_quantiles=32)

cql = CQL(q_func_factory=q_func)
```

The default Q function is mean approximator, which estimates expected scalar action-values. However, in recent advancements of deep reinforcement learning, the new type of action-value approximators has been proposed, which is called *distributional* Q functions.

Unlike the `mean` approximator, the *distributional* Q functions estimate distribution of action-values. This *distributional* approaches have shown consistently much stronger performance than the `mean` approximator.

Here is a list of available Q functions in the order of performance ascendingly. Currently, as a trade-off between performance and computational complexity, the higher performance requires the more expensive computational costs.

<code>d3rlpy.models.q_functions.</code> <code>MeanQFunctionFactory</code>	Standard Q function factory class.
<code>d3rlpy.models.q_functions.</code> <code>QRQFunctionFactory</code>	Quantile Regression Q function factory class.
<code>d3rlpy.models.q_functions.</code> <code>IQNQFunctionFactory</code>	Implicit Quantile Network Q function factory class.
<code>d3rlpy.models.q_functions.</code> <code>FQFQFunctionFactory</code>	Fully parameterized Quantile Function Q function factory.

3.2.1 d3rlpy.models.q_functions.MeanQFunctionFactory

class d3rlpy.models.q_functions.**MeanQFunctionFactory** (*bootstrap=False*,
share_encoder=False)

Standard Q function factory class.

This is the standard Q function factory class.

References

- Mnih et al., Human-level control through deep reinforcement learning.
- Lillicrap et al., Continuous control with deep reinforcement learning.

Parameters

- **bootstrap** (*bool*) – flag to bootstrap Q functions.
- **share_encoder** (*bool*) – flag to share encoder over multiple Q functions.

Methods

create_continuous (*encoder*)

Returns PyTorch's Q function module.

Parameters **encoder** (*d3rlpy.models.torch.encoders.EncoderWithAction*)
– an encoder module that processes the observation and action to obtain feature representations.

Returns continuous Q function object.

Return type d3rlpy.models.torch.q_functions.ContinuousMeanQFunction

create_discrete (*encoder, action_size*)

Returns PyTorch's Q function module.

Parameters

- **encoder** (*d3rlpy.models.torch.encoders.Encoder*) – an encoder module that processes the observation to obtain feature representations.
- **action_size** (*int*) – dimension of discrete action-space.

Returns discrete Q function object.

Return type d3rlpy.models.torch.q_functions.DiscreteMeanQFunction

get_params (*deep=False*)

Returns Q function parameters.

Returns Q function parameters.

Parameters **deep** (*bool*) –

Return type Dict[str, Any]

get_type ()

Returns Q function type.

Returns Q function type.

Return type str

Attributes

TYPE: `ClassVar[str] = 'mean'`

bootstrap

share_encoder

3.2.2 d3rlpy.models.q_functions.QRQFunctionFactory

```
class d3rlpy.models.q_functions.QRQFunctionFactory(bootstrap=False,  
                                                    share_encoder=False,  
                                                    n_quantiles=32)
```

Quantile Regression Q function factory class.

References

- Dabney et al., Distributional reinforcement learning with quantile regression.

Parameters

- **bootstrap** (*bool*) – flag to bootstrap Q functions.
- **share_encoder** (*bool*) – flag to share encoder over multiple Q functions.
- **n_quantiles** – the number of quantiles.

Methods

create_continuous (*encoder*)

Returns PyTorch's Q function module.

Parameters **encoder** (*d3rlpy.models.torch.encoders.EncoderWithAction*)
– an encoder module that processes the observation and action to obtain feature representations.

Returns continuous Q function object.

Return type *d3rlpy.models.torch.q_functions.ContinuousQRQFunction*

create_discrete (*encoder, action_size*)

Returns PyTorch's Q function module.

Parameters

- **encoder** (*d3rlpy.models.torch.encoders.Encoder*) – an encoder module that processes the observation to obtain feature representations.
- **action_size** (*int*) – dimension of discrete action-space.

Returns discrete Q function object.

Return type *d3rlpy.models.torch.q_functions.DiscreteQRQFunction*

get_params (*deep=False*)

Returns Q function parameters.

Returns Q function parameters.

Parameters **deep** (*bool*) –

Return type Dict[str, Any]

get_type()

Returns Q function type.

Returns Q function type.

Return type str

Attributes

TYPE: ClassVar[str] = 'qr'

bootstrap

n_quantiles

share_encoder

3.2.3 d3rlpy.models.q_functions.IQNQFunctionFactory

```
class d3rlpy.models.q_functions.IQNQFunctionFactory (bootstrap=False,
                                                    share_encoder=False,
                                                    n_quantiles=64,
                                                    n_greedy_quantiles=32,    em-
                                                    bed_size=64)
```

Implicit Quantile Network Q function factory class.

References

- Dabney et al., Implicit quantile networks for distributional reinforcement learning.

Parameters

- **bootstrap** (bool) – flag to bootstrap Q functions.
- **share_encoder** (bool) – flag to share encoder over multiple Q functions.
- **n_quantiles** – the number of quantiles.
- **n_greedy_quantiles** – the number of quantiles for inference.
- **embed_size** – the embedding size.

Methods

create_continuous (encoder)

Returns PyTorch's Q function module.

Parameters **encoder** (d3rlpy.models.torch.encoders.EncoderWithAction)
– an encoder module that processes the observation and action to obtain feature representations.

Returns continuous Q function object.

Return type d3rlpy.models.torch.q_functions.ContinuousIQNQFunction

create_discrete (encoder, action_size)

Returns PyTorch's Q function module.

Parameters

- **encoder** (*d3rlpy.models.torch.encoders.Encoder*) – an encoder module that processes the observation to obtain feature representations.
- **action_size** (*int*) – dimension of discrete action-space.

Returns discrete Q function object.

Return type *d3rlpy.models.torch.q_functions.DiscreteIQNQFunction*

get_params (*deep=False*)

Returns Q function parameters.

Returns Q function parameters.

Parameters **deep** (*bool*) –

Return type *Dict[str, Any]*

get_type ()

Returns Q function type.

Returns Q function type.

Return type *str*

Attributes

TYPE: *ClassVar[str]* = 'iqn'

bootstrap

embed_size

n_greedy_quantiles

n_quantiles

share_encoder

3.2.4 d3rlpy.models.q_functions.FQFQFunctionFactory

```
class d3rlpy.models.q_functions.FQFQFunctionFactory(bootstrap=False,
                                                    share_encoder=False,
                                                    n_quantiles=32,
                                                    embed_size=64,
                                                    entropy_coeff=0.0)
```

Fully parameterized Quantile Function Q function factory.

References

- Yang et al., Fully parameterized quantile function for distributional reinforcement learning.

Parameters

- **bootstrap** (*bool*) – flag to bootstrap Q functions.
- **share_encoder** (*bool*) – flag to share encoder over multiple Q functions.
- **n_quantiles** – the number of quantiles.

- **embed_size** – the embedding size.
- **entropy_coeff** – the coefficient of entropy penalty term.

Methods

create_continuous (*encoder*)

Returns PyTorch's Q function module.

Parameters **encoder** (*d3rlpy.models.torch.encoders.EncoderWithAction*)
– an encoder module that processes the observation and action to obtain feature representations.

Returns continuous Q function object.

Return type *d3rlpy.models.torch.q_functions.ContinuousFQFQFunction*

create_discrete (*encoder, action_size*)

Returns PyTorch's Q function module.

Parameters

- **encoder** (*d3rlpy.models.torch.encoders.Encoder*) – an encoder module that processes the observation to obtain feature representations.
- **action_size** (*int*) – dimension of discrete action-space.

Returns discrete Q function object.

Return type *d3rlpy.models.torch.q_functions.DiscreteFQFQFunction*

get_params (*deep=False*)

Returns Q function parameters.

Returns Q function parameters.

Parameters **deep** (*bool*) –

Return type *Dict[str, Any]*

get_type ()

Returns Q function type.

Returns Q function type.

Return type *str*

Attributes

TYPE: *ClassVar[str]* = 'fqf'

bootstrap

embed_size

entropy_coeff

n_quantiles

share_encoder

3.3 MDPDataset

d3rlpy provides useful dataset structure for data-driven deep reinforcement learning. In supervised learning, the training script iterates input data X and label data Y . However, in reinforcement learning, mini-batches consist with sets of $(s_t, a_t, r_{t+1}, s_{t+1})$ and episode terminal flags. Converting a set of observations, actions, rewards and terminal flags into this tuples is boring and requires some codings.

Therefore, d3rlpy provides *MDPDataset* class which enables you to handle reinforcement learning datasets without any efforts.

```
from d3rlpy.dataset import MDPDataset

# 1000 steps of observations with shape of (100,)
observations = np.random.random((1000, 100))
# 1000 steps of actions with shape of (4,)
actions = np.random.random((1000, 4))
# 1000 steps of rewards
rewards = np.random.random(1000)
# 1000 steps of terminal flags
terminals = np.random.randint(2, size=1000)

dataset = MDPDataset(observations, actions, rewards, terminals)

# automatically splitted into d3rlpy.dataset.Episode objects
dataset.episodes

# each episode is also splitted into d3rlpy.dataset.Transition objects
episode = dataset.episodes[0]
episode[0].observation
episode[0].action
episode[0].next_reward
episode[0].next_observation
episode[0].terminal

# d3rlpy.dataset.Transition object has pointers to previous and next
# transitions like linked list.
transition = episode[0]
while transition.next_transition:
    transition = transition.next_transition

# save as HDF5
dataset.dump('dataset.h5')

# load from HDF5
new_dataset = MDPDataset.load('dataset.h5')
```

<i>d3rlpy.dataset.MDPDataset</i>	Markov-Decision Process Dataset class.
<i>d3rlpy.dataset.Episode</i>	Episode class.
<i>d3rlpy.dataset.Transition</i>	Transition class.
<i>d3rlpy.dataset.TransitionMiniBatch</i>	mini-batch of Transition objects.

3.3.1 d3rlpy.dataset.MDPDataset

class d3rlpy.dataset.MDPDataset (observations, actions, rewards, terminals, episode_terminals=None, discrete_action=None, create_mask=False, mask_size=1)

Markov-Decision Process Dataset class.

MDPDataset is designed for reinforcement learning datasets to use them like supervised learning datasets.

```
from d3rlpy.dataset import MDPDataset

# 1000 steps of observations with shape of (100,)
observations = np.random.random((1000, 100))
# 1000 steps of actions with shape of (4,)
actions = np.random.random((1000, 4))
# 1000 steps of rewards
rewards = np.random.random(1000)
# 1000 steps of terminal flags
terminals = np.random.randint(2, size=1000)

dataset = MDPDataset(observations, actions, rewards, terminals)
```

The MDPDataset object automatically splits the given data into list of `d3rlpy.dataset.Episode` objects. Furthermore, the MDPDataset object behaves like a list in order to use with scikit-learn utilities.

```
# returns the number of episodes
len(dataset)

# access to the first episode
episode = dataset[0]

# iterate through all episodes
for episode in dataset:
    pass
```

Parameters

- **observations** (`numpy.ndarray`) – N-D array. If the observation is a vector, the shape should be $(N, \text{dim_observation})$. If the observations is an image, the shape should be (N, C, H, W) .
- **actions** (`numpy.ndarray`) – N-D array. If the actions-space is continuous, the shape should be $(N, \text{dim_action})$. If the action-space is discrete, the shape should be $(N,)$.
- **rewards** (`numpy.ndarray`) – array of scalar rewards.
- **terminals** (`numpy.ndarray`) – array of binary terminal flags.
- **episode_terminals** (`numpy.ndarray`) – array of binary episode terminal flags. The given data will be splitted based on this flag. This is useful if you want to specify the non-environment terminations (e.g. timeout). If `None`, the episode terminations match the environment terminations.
- **discrete_action** (`bool`) – flag to use the given actions as discrete action-space actions. If `None`, the action type is automatically determined.
- **create_mask** (`bool`) – flag to create binary masks for bootstrapping.
- **mask_size** (`int`) – ensemble size for mask. If `create_mask` is `False`, this will be ignored.

Methods

`__getitem__` (*index*)

`__len__` ()

`__iter__` ()

append (*observations, actions, rewards, terminals, episode_terminals=None*)

Appends new data.

Parameters

- **observations** (*numpy.ndarray*) – N-D array.
- **actions** (*numpy.ndarray*) – actions.
- **rewards** (*numpy.ndarray*) – rewards.
- **terminals** (*numpy.ndarray*) – terminals.
- **episode_terminals** (*numpy.ndarray*) – episode terminals.

build_episodes ()

Builds episode objects.

This method will be internally called when accessing the episodes property at the first time.

clip_reward (*low=None, high=None*)

Clips rewards in the given range.

Parameters

- **low** (*float*) – minimum value. If None, clipping is not performed on lower edge.
- **high** (*float*) – maximum value. If None, clipping is not performed on upper edge.

compute_stats ()

Computes statistics of the dataset.

```
stats = dataset.compute_stats()

# return statistics
stats['return']['mean']
stats['return']['std']
stats['return']['min']
stats['return']['max']

# reward statistics
stats['reward']['mean']
stats['reward']['std']
stats['reward']['min']
stats['reward']['max']

# action (only with continuous control actions)
stats['action']['mean']
stats['action']['std']
stats['action']['min']
stats['action']['max']

# observation (only with numpy.ndarray observations)
stats['observation']['mean']
stats['observation']['std']
```

(continues on next page)

(continued from previous page)

```
stats['observation']['min']
stats['observation']['max']
```

Returns statistics of the dataset.

Return type `dict`

dump (*fname*)

Saves dataset as HDF5.

Parameters **fname** (*str*) – file path.

extend (*dataset*)

Extend dataset by another dataset.

Parameters **dataset** (`d3rlpy.dataset.MDPDataset`) – dataset.

get_action_size ()

Returns dimension of action-space.

If *discrete_action=True*, the return value will be the maximum index +1 in the give actions.

Returns dimension of action-space.

Return type `int`

get_observation_shape ()

Returns observation shape.

Returns observation shape.

Return type `tuple`

is_action_discrete ()

Returns *discrete_action* flag.

Returns *discrete_action* flag.

Return type `bool`

classmethod load (*fname*, *create_mask=False*, *mask_size=1*)

Loads dataset from HDF5.

```
import numpy as np
from d3rlpy.dataset import MDPDataset

dataset = MDPDataset(np.random.random(10, 4),
                     np.random.random(10, 2),
                     np.random.random(10),
                     np.random.randint(2, size=10))

# save as HDF5
dataset.dump('dataset.h5')

# load from HDF5
new_dataset = MDPDataset.load('dataset.h5')
```

Parameters

- **fname** (*str*) – file path.
- **create_mask** (*bool*) – flag to create bootstrapping masks.

- **mask_size** (*int*) – size of bootstrapping masks.

size()

Returns the number of episodes in the dataset.

Returns the number of episodes.

Return type *int*

Attributes

actions

Returns the actions.

Returns array of actions.

Return type *numpy.ndarray*

episode_terminals

Returns the episode terminal flags.

Returns array of episode terminal flags.

Return type *numpy.ndarray*

episodes

Returns the episodes.

Returns list of *d3rlpy.dataset.Episode* objects.

Return type *list(d3rlpy.dataset.Episode)*

observations

Returns the observations.

Returns array of observations.

Return type *numpy.ndarray*

rewards

Returns the rewards.

Returns array of rewards

Return type *numpy.ndarray*

terminals

Returns the terminal flags.

Returns array of terminal flags.

Return type *numpy.ndarray*

3.3.2 d3rlpy.dataset.Episode

class d3rlpy.dataset.**Episode** (*observation_shape, action_size, observations, actions, rewards, terminal=True, create_mask=False, mask_size=1*)

Episode class.

This class is designed to hold data collected in a single episode.

Episode object automatically splits data into list of `d3rlpy.dataset.Transition` objects. Also Episode object behaves like a list object for ease of access to transitions.

```
# return the number of transitions
len(episode)

# access to the first transition
transitions = episode[0]

# iterate through all transitions
for transition in episode:
    pass
```

Parameters

- **observation_shape** (*tuple*) – observation shape.
- **action_size** (*int*) – dimension of action-space.
- **observations** (*numpy.ndarray*) – observations.
- **actions** (*numpy.ndarray*) – actions.
- **rewards** (*numpy.ndarray*) – scalar rewards.
- **terminal** (*bool*) – binary terminal flag. If False, the episode is not terminated by the environment (e.g. timeout).
- **create_mask** (*bool*) – flag to create binary masks for bootstrapping.
- **mask_size** (*int*) – ensemble size for mask. If `create_mask` is False, this will be ignored.

Methods

__getitem__ (*index*)

__len__ ()

__iter__ ()

build_transitions ()

Builds transition objects.

This method will be internally called when accessing the transitions property at the first time.

compute_return ()

Computes sum of rewards.

$$R = \sum_{i=1} r_i$$

Returns episode return.

Return type `float`

get_action_size()

Returns dimension of action-space.

Returns dimension of action-space.

Return type `int`

get_observation_shape()

Returns observation shape.

Returns observation shape.

Return type `tuple`

size()

Returns the number of transitions.

Returns the number of transitions.

Return type `int`

Attributes

actions

Returns the actions.

Returns array of actions.

Return type `numpy.ndarray`

observations

Returns the observations.

Returns array of observations.

Return type `numpy.ndarray`

rewards

Returns the rewards.

Returns array of rewards.

Return type `numpy.ndarray`

terminal

Returns the terminal flag.

Returns the terminal flag.

Return type `bool`

transitions

Returns the transitions.

Returns list of `d3rlpy.dataset.Transition` objects.

Return type `list(d3rlpy.dataset.Transition)`

3.3.3 d3rlpy.dataset.Transition

class d3rlpy.dataset.Transition

Transition class.

This class is designed to hold data between two time steps, which is usually used as inputs of loss calculation in reinforcement learning.

Parameters

- **observation_shape** (*tuple*) – observation shape.
- **action_size** (*int*) – dimension of action-space.
- **observation** (*numpy.ndarray*) – observation at t .
- **action** (*numpy.ndarray* or *int*) – action at t .
- **reward** (*float*) – reward at t .
- **next_observation** (*numpy.ndarray*) – observation at $t+1$.
- **next_action** (*numpy.ndarray* or *int*) – action at $t+1$.
- **next_reward** (*float*) – reward at $t+1$.
- **terminal** (*int*) – terminal flag at $t+1$.
- **mask** (*numpy.ndarray*) – binary mask for bootstrapping.
- **prev_transition** (*d3rlpy.dataset.Transition*) – pointer to the previous transition.
- **next_transition** (*d3rlpy.dataset.Transition*) – pointer to the next transition.

Methods

clear_links ()

Clears links to the next and previous transitions.

This method is necessary to call when freeing this instance by GC.

get_action_size ()

Returns dimension of action-space.

Returns dimension of action-space.

Return type *int*

get_observation_shape ()

Returns observation shape.

Returns observation shape.

Return type *tuple*

Attributes

action

Returns action at t .

Returns action at t .

Return type (`numpy.ndarray` or `int`)

mask

Returns binary mask for bootstrapping.

Returns array of binary mask.

Return type `np.ndarray`

next_action

Returns action at $t+1$.

Returns action at $t+1$.

Return type (`numpy.ndarray` or `int`)

next_observation

Returns observation at $t+1$.

Returns observation at $t+1$.

Return type `numpy.ndarray` or `torch.Tensor`

next_reward

Returns reward at $t+1$.

Returns reward at $t+1$.

Return type `float`

next_transition

Returns pointer to the next transition.

If this is the last transition, this method should return `None`.

Returns next transition.

Return type `d3rlpy.dataset.Transition`

observation

Returns observation at t .

Returns observation at t .

Return type `numpy.ndarray` or `torch.Tensor`

prev_transition

Returns pointer to the previous transition.

If this is the first transition, this method should return `None`.

Returns previous transition.

Return type `d3rlpy.dataset.Transition`

reward

Returns reward at t .

Returns reward at t .

Return type `float`

terminal

Returns terminal flag at $t+1$.

Returns terminal flag at $t+1$.

Return type `int`

3.3.4 d3rlpy.dataset.TransitionMiniBatch

class `d3rlpy.dataset.TransitionMiniBatch`

mini-batch of Transition objects.

This class is designed to hold `d3rlpy.dataset.Transition` objects for being passed to algorithms during fitting.

If the observation is image, you can stack arbitrary frames via `n_frames`.

```
transition.observation.shape == (3, 84, 84)

batch_size = len(transitions)

# stack 4 frames
batch = TransitionMiniBatch(transitions, n_frames=4)

# 4 frames x 3 channels
batch.observations.shape == (batch_size, 12, 84, 84)
```

This is implemented by tracing previous transitions through `prev_transition` property.

Parameters

- **transitions** (`list(d3rlpy.dataset.Transition)`) – mini-batch of transitions.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – length of N-step sampling.
- **gamma** (`float`) – discount factor for N-step calculation.

Methods

`__getitem__` (`key`, /)
Return `self[key]`.

`__len__` ()
Return `len(self)`.

`__iter__` ()
Implement `iter(self)`.

`add_additional_data` (`key`, `value`)
Add arbitrary additional data.

Parameters

- **key** (`str`) – key of data.
- **value** (`any`) – value.

get_additional_data (*key*)

Returns specified additional data.

Parameters **key** (*str*) – key of data.

Returns value.

Return type any

size ()

Returns size of mini-batch.

Returns mini-batch size.

Return type int

Attributes

actions

Returns mini-batch of actions at t .

Returns actions at t .

Return type `numpy.ndarray`

masks

Returns mini-batch of binary masks for bootstrapping.

If any of transitions have an invalid mask, this will return `None`.

Returns binary mask.

Return type `numpy.ndarray`

n_steps

Returns mini-batch of the number of steps before next observations.

This will always include only ones if `n_steps=1`. If `n_steps` is bigger than 1. the values will depend on its episode length.

Returns the number of steps before next observations.

Return type `numpy.ndarray`

next_actions

Returns mini-batch of actions at $t+n$.

Returns actions at $t+n$.

Return type `numpy.ndarray`

next_observations

Returns mini-batch of observations at $t+n$.

Returns observations at $t+n$.

Return type `numpy.ndarray` or `torch.Tensor`

next_rewards

Returns mini-batch of rewards at $t+n$.

Returns rewards at $t+n$.

Return type `numpy.ndarray`

observations

Returns mini-batch of observations at t .

Returns observations at t .

Return type `numpy.ndarray` or `torch.Tensor`

rewards

Returns mini-batch of rewards at t .

Returns rewards at t .

Return type `numpy.ndarray`

terminals

Returns mini-batch of terminal flags at $t+n$.

Returns terminal flags at $t+n$.

Return type `numpy.ndarray`

transitions

Returns transitions.

Returns list of transitions.

Return type `d3rlpy.dataset.Transition`

3.4 Datasets

d3rlpy provides datasets for experimenting data-driven deep reinforcement learning algorithms.

<code>d3rlpy.datasets.get_cartpole</code>	Returns cartpole dataset and environment.
<code>d3rlpy.datasets.get_pendulum</code>	Returns pendulum dataset and environment.
<code>d3rlpy.datasets.get_pybullet</code>	Returns pybullet dataset and environment.
<code>d3rlpy.datasets.get_atari</code>	Returns atari dataset and environment.
<code>d3rlpy.datasets.get_d4rl</code>	Returns d4rl dataset and environment.
<code>d3rlpy.datasets.get_dataset</code>	Returns dataset and environment by guessing from name.

3.4.1 d3rlpy.datasets.get_cartpole

`d3rlpy.datasets.get_cartpole` (*create_mask=False*, *mask_size=1*, *dataset_type='replay'*)

Returns cartpole dataset and environment.

The dataset is automatically downloaded to `d3rlpy_data/cartpole.h5` if it does not exist.

Parameters

- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.
- **dataset_type** (*str*) – dataset type. Available options are `['replay', 'random']`.

Returns tuple of `d3rlpy.dataset.MDPDataset` and gym environment.

Return type `Tuple[d3rlpy.dataset.MDPDataset, gym.core.Env]`

3.4.2 d3rlpy.datasets.get_pendulum

`d3rlpy.datasets.get_pendulum(create_mask=False, mask_size=1, dataset_type='replay')`

Returns pendulum dataset and environment.

The dataset is automatically downloaded to `d3rlpy_data/pendulum.h5` if it does not exist.

Parameters

- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.
- **dataset_type** (*str*) – dataset type. Available options are `['replay', 'random']`.

Returns tuple of `d3rlpy.dataset.MDPDataset` and gym environment.

Return type `Tuple[d3rlpy.dataset.MDPDataset, gym.core.Env]`

3.4.3 d3rlpy.datasets.get_pybullet

`d3rlpy.datasets.get_pybullet(env_name, create_mask=False, mask_size=1)`

Returns pybullet dataset and environment.

The dataset is provided through `d4rl-pybullet`. See more details including available dataset from its GitHub page.

```
from d3rlpy.datasets import get_pybullet

dataset, env = get_pybullet('hopper-bullet-mixed-v0')
```

References

- <https://github.com/takuseno/d4rl-pybullet>

Parameters

- **env_name** (*str*) – environment id of `d4rl-pybullet` dataset.
- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.

Returns tuple of `d3rlpy.dataset.MDPDataset` and gym environment.

Return type `Tuple[d3rlpy.dataset.MDPDataset, gym.core.Env]`

3.4.4 d3rlpy.datasets.get_atari

`d3rlpy.datasets.get_atari(env_name, create_mask=False, mask_size=1)`

Returns atari dataset and environment.

The dataset is provided through `d4rl-atari`. See more details including available dataset from its GitHub page.

```
from d3rlpy.datasets import get_atari

dataset, env = get_atari('breakout-mixed-v0')
```

References

- <https://github.com/takuseno/d4rl-atari>

Parameters

- **env_name** (*str*) – environment id of d4rl-atari dataset.
- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.

Returns tuple of *d3rlpy.dataset.MDPDataset* and gym environment.

Return type Tuple[*d3rlpy.dataset.MDPDataset*, gym.core.Env]

3.4.5 d3rlpy.datasets.get_d4rl

d3rlpy.datasets.get_d4rl (*env_name*, *create_mask=False*, *mask_size=1*)

Returns d4rl dataset and environment.

The dataset is provided through d4rl.

```
from d3rlpy.datasets import get_d4rl

dataset, env = get_d4rl('hopper-medium-v0')
```

References

- Fu et al., D4RL: Datasets for Deep Data-Driven Reinforcement Learning.
- <https://github.com/rail-berkeley/d4rl>

Parameters

- **env_name** (*str*) – environment id of d4rl dataset.
- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.

Returns tuple of *d3rlpy.dataset.MDPDataset* and gym environment.

Return type Tuple[*d3rlpy.dataset.MDPDataset*, gym.core.Env]

3.4.6 d3rlpy.datasets.get_dataset

d3rlpy.datasets.get_dataset (*env_name*, *create_mask=False*, *mask_size=1*)

Returns dataset and environment by guessing from name.

This function returns dataset by matching name with the following datasets.

- cartpole
- pendulum
- d4rl-pybullet
- d4rl-atari

- d4rl

```
import d3rlpy

# cartpole dataset
dataset, env = d3rlpy.datasets.get_dataset('cartpole')

# pendulum dataset
dataset, env = d3rlpy.datasets.get_dataset('pendulum')

# d4rl-pybullet dataset
dataset, env = d3rlpy.datasets.get_dataset('hopper-bullet-mixed-v0')

# d4rl-atari dataset
dataset, env = d3rlpy.datasets.get_dataset('breakout-mixed-v0')

# d4rl dataset
dataset, env = d3rlpy.datasets.get_dataset('hopper-medium-v0')
```

Parameters

- **env_name** (*str*) – environment id of the dataset.
- **create_mask** (*bool*) – flag to create binary mask for bootstrapping.
- **mask_size** (*int*) – ensemble size for binary mask.

Returns tuple of `d3rlpy.dataset.MDPDataset` and gym environment.

Return type Tuple[`d3rlpy.dataset.MDPDataset`, `gym.core.Env`]

3.5 Preprocessing

3.5.1 Observation

d3rlpy provides several preprocessors tightly incorporated with algorithms. Each preprocessor is implemented with PyTorch operation, which will be included in the model exported by `save_policy` method.

```
from d3rlpy.algos import CQL
from d3rlpy.dataset import MDPDataset

dataset = MDPDataset(...)

# choose from ['pixel', 'min_max', 'standard'] or None
cql = CQL(scaler='standard')

# scaler is fitted from the given episodes
cql.fit(dataset.episodes)

# preprocessing is included in TorchScript
cql.save_policy('policy.pt')

# you don't need to take care of preprocessing at production
policy = torch.jit.load('policy.pt')
action = policy(unpreprocessed_x)
```

You can also initialize scalers by yourself.

```
from d3rlpy.preprocessing import StandardScaler

scaler = StandardScaler(mean=..., std=...)

cql = CQL(scaler=scaler)
```

<code>d3rlpy.preprocessing.PixelScaler</code>	Pixel normalization preprocessing.
<code>d3rlpy.preprocessing.MinMaxScaler</code>	Min-Max normalization preprocessing.
<code>d3rlpy.preprocessing.StandardScaler</code>	Standardization preprocessing.

d3rlpy.preprocessing.PixelScaler

class d3rlpy.preprocessing.**PixelScaler**
Pixel normalization preprocessing.

$$x' = x/255$$

```
from d3rlpy.dataset import MDPDataset
from d3rlpy.algos import CQL

dataset = MDPDataset(observations, actions, rewards, terminals)

# initialize algorithm with PixelScaler
cql = CQL(scaler='pixel')

cql.fit(dataset.episodes)
```

Methods

fit (*episodes*)

Estimates scaling parameters from dataset.

Parameters **episodes** (*List* [`d3rlpy.dataset.Episode`]) – list of episodes.

Return type `None`

fit_with_env (*env*)

Gets scaling parameters from environment.

Parameters **env** (*`gym.core.Env`*) – gym environment.

Return type `None`

get_params (*deep=False*)

Returns scaling parameters.

Parameters **deep** (*`bool`*) – flag to deeply copy objects.

Returns scaler parameters.

Return type `Dict[str, Any]`

get_type ()

Returns a scaler type.

Returns scaler type.

Return type `str`

reverse_transform (x)

Returns reversely transformed observations.

Parameters \mathbf{x} (`torch.Tensor`) – observation.

Returns reversely transformed observation.

Return type `torch.Tensor`

transform (x)

Returns processed observations.

Parameters \mathbf{x} (`torch.Tensor`) – observation.

Returns processed observation.

Return type `torch.Tensor`

Attributes

TYPE: `ClassVar[str] = 'pixel'`

d3rlpy.preprocessing.MinMaxScaler

class `d3rlpy.preprocessing.MinMaxScaler` (*dataset=None, maximum=None, minimum=None*)
Min-Max normalization preprocessing.

$$x' = (x - \min x) / (\max x - \min x)$$

```
from d3rlpy.dataset import MDPDataset
from d3rlpy.algos import CQL

dataset = MDPDataset(observations, actions, rewards, terminals)

# initialize algorithm with MinMaxScaler
cql = CQL(scaler='min_max')

# scaler is initialized from the given episodes
cql.fit(dataset.episodes)
```

You can also initialize with `d3rlpy.dataset.MDPDataset` object or manually.

```
from d3rlpy.preprocessing import MinMaxScaler

# initialize with dataset
scaler = MinMaxScaler(dataset)

# initialize manually
minimum = observations.min(axis=0)
maximum = observations.max(axis=0)
scaler = MinMaxScaler(minimum=minimum, maximum=maximum)

cql = CQL(scaler=scaler)
```

Parameters

- **dataset** (`d3rlpy.dataset.MDPDataset`) – dataset object.
- **min** (`numpy.ndarray`) – minimum values at each entry.
- **max** (`numpy.ndarray`) – maximum values at each entry.

Methods

fit (*episodes*)

Estimates scaling parameters from dataset.

Parameters **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Return type `None`

fit_with_env (*env*)

Gets scaling parameters from environment.

Parameters **env** (`gym.core.Env`) – gym environment.

Return type `None`

get_params (*deep=False*)

Returns scaling parameters.

Parameters **deep** (`bool`) – flag to deeply copy objects.

Returns scaler parameters.

Return type `Dict[str, Any]`

get_type ()

Returns a scaler type.

Returns scaler type.

Return type `str`

reverse_transform (*x*)

Returns reversely transformed observations.

Parameters **x** (`torch.Tensor`) – observation.

Returns reversely transformed observation.

Return type `torch.Tensor`

transform (*x*)

Returns processed observations.

Parameters **x** (`torch.Tensor`) – observation.

Returns processed observation.

Return type `torch.Tensor`

Attributes

TYPE: `ClassVar[str] = 'min_max'`

d3rlpy.preprocessing.StandardScaler

class d3rlpy.preprocessing.**StandardScaler** (*dataset=None, mean=None, std=None*)
Standardization preprocessing.

$$x' = (x - \mu) / \sigma$$

```
from d3rlpy.dataset import MDPDataset
from d3rlpy.algos import CQL

dataset = MDPDataset(observations, actions, rewards, terminals)

# initialize algorithm with StandardScaler
cql = CQL(scaler='standard')

# scaler is initialized from the given episodes
cql.fit(dataset.episodes)
```

You can initialize with `d3rlpy.dataset.MDPDataset` object or manually.

```
from d3rlpy.preprocessing import StandardScaler

# initialize with dataset
scaler = StandardScaler(dataset)

# initialize manually
mean = observations.mean(axis=0)
std = observations.std(axis=0)
scaler = StandardScaler(mean=mean, std=std)

cql = CQL(scaler=scaler)
```

Parameters

- **dataset** (`d3rlpy.dataset.MDPDataset`) – dataset object.
- **mean** (`numpy.ndarray`) – mean values at each entry.
- **std** (`numpy.ndarray`) – standard deviation at each entry.

Methods

fit (*episodes*)

Estimates scaling parameters from dataset.

Parameters **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Return type `None`

fit_with_env (*env*)

Gets scaling parameters from environment.

Parameters **env** (`gym.core.Env`) – gym environment.

Return type `None`

get_params (*deep=False*)
Returns scaling parameters.

Parameters *deep* (*bool*) – flag to deeply copy objects.

Returns scaler parameters.

Return type `Dict[str, Any]`

get_type ()
Returns a scaler type.

Returns scaler type.

Return type `str`

reverse_transform (*x*)
Returns reversely transformed observations.

Parameters *x* (*torch.Tensor*) – observation.

Returns reversely transformed observation.

Return type `torch.Tensor`

transform (*x*)
Returns processed observations.

Parameters *x* (*torch.Tensor*) – observation.

Returns processed observation.

Return type `torch.Tensor`

Attributes

TYPE: `ClassVar[str] = 'standard'`

3.5.2 Action

d3rlpy also provides the feature that preprocesses continuous action. With this preprocessing, you don't need to normalize actions in advance or implement normalization in the environment side.

```
from d3rlpy.algos import CQL
from d3rlpy.dataset import MDPDataset

dataset = MDPDataset(...)

# 'min_max' or None
cql = CQL(action_scaler='min_max')

# action scaler is fitted from the given episodes
cql.fit(dataset.episodes)

# postprocessing is included in TorchScript
cql.save_policy('policy.pt')

# you don't need to take care of postprocessing at production
```

(continues on next page)

(continued from previous page)

```
policy = torch.jit.load('policy.pt')
action = policy(x)
```

You can also initialize scalers by yourself.

```
from d3rlpy.preprocessing import MinMaxActionScaler

action_scaler = MinMaxActionScaler(minimum=..., maximum=...)

cql = CQL(action_scaler=action_scaler)
```

d3rlpy.preprocessing.
MinMaxActionScaler

Min-Max normalization action preprocessing.

d3rlpy.preprocessing.MinMaxActionScaler

class d3rlpy.preprocessing.MinMaxActionScaler(*dataset=None, maximum=None, minimum=None*)

Min-Max normalization action preprocessing.

Actions will be normalized in range $[-1.0, 1.0]$.

$$a' = (a - \min a) / (\max a - \min a) * 2 - 1$$

```
from d3rlpy.dataset import MDPDataset
from d3rlpy.algos import CQL

dataset = MDPDataset(observations, actions, rewards, terminals)

# initialize algorithm with MinMaxActionScaler
cql = CQL(action_scaler='min_max')

# scaler is initialized from the given episodes
cql.fit(dataset.episodes)
```

You can also initialize with *d3rlpy.dataset.MDPDataset* object or manually.

```
from d3rlpy.preprocessing import MinMaxActionScaler

# initialize with dataset
scaler = MinMaxActionScaler(dataset)

# initialize manually
minimum = actions.min(axis=0)
maximum = actions.max(axis=0)
action_scaler = MinMaxActionScaler(minimum=minimum, maximum=maximum)

cql = CQL(action_scaler=action_scaler)
```

Parameters

- **dataset** (*d3rlpy.dataset.MDPDataset*) – dataset object.
- **min** (*numpy.ndarray*) – minimum values at each entry.

- **max** (*numpy.ndarray*) – maximum values at each entry.

Methods

fit (*episodes*)

Estimates scaling parameters from dataset.

Parameters **episodes** (*List [d3rlpy.dataset.Episode]*) – a list of episode objects.

Return type *None*

fit_with_env (*env*)

Gets scaling parameters from environment.

Parameters **env** (*gym.core.Env*) – gym environment.

Return type *None*

get_params (*deep=False*)

Returns action scaler params.

Parameters **deep** (*bool*) – flag to deepcopy parameters.

Returns action scaler parameters.

Return type *Dict[str, Any]*

get_type ()

Returns action scaler type.

Returns action scaler type.

Return type *str*

reverse_transform (*action*)

Returns reversely transformed action.

Parameters **action** (*torch.Tensor*) – action vector.

Returns reversely transformed action.

Return type *torch.Tensor*

reverse_transform_numpy (*action*)

Returns reversely transformed action in numpy array.

Parameters **action** (*numpy.ndarray*) – action vector.

Returns reversely transformed action.

Return type *numpy.ndarray*

transform (*action*)

Returns processed action.

Parameters **action** (*torch.Tensor*) – action vector.

Returns processed action.

Return type *torch.Tensor*

Attributes

TYPE: `ClassVar[str]` = 'min_max'

3.6 Optimizers

d3rlpy provides `OptimizerFactory` that gives you flexible control over optimizers. `OptimizerFactory` takes PyTorch's optimizer class and its arguments to initialize, which you can check more [here](#).

```
from torch.optim import Adam
from d3rlpy.algos import DQN
from d3rlpy.models.optimizers import OptimizerFactory

# modify weight decay
optim_factory = OptimizerFactory(Adam, weight_decay=1e-4)

# set OptimizerFactory
dqn = DQN(optim_factory=optim_factory)
```

There are also convenient aliases.

```
from d3rlpy.models.optimizers import AdamFactory

# alias for Adam optimizer
optim_factory = AdamFactory(weight_decay=1e-4)

dqn = DQN(optim_factory=optim_factory)
```

<code>d3rlpy.models.optimizers.OptimizerFactory</code>	A factory class that creates an optimizer object in a lazy way.
<code>d3rlpy.models.optimizers.SGDFactory</code>	An alias for SGD optimizer.
<code>d3rlpy.models.optimizers.AdamFactory</code>	An alias for Adam optimizer.
<code>d3rlpy.models.optimizers.RMSpropFactory</code>	An alias for RMSprop optimizer.

3.6.1 d3rlpy.models.optimizers.OptimizerFactory

class `d3rlpy.models.optimizers.OptimizerFactory` (*optim_cls, **kwargs*)

A factory class that creates an optimizer object in a lazy way.

The optimizers in algorithms can be configured through this factory class.

```
from torch.optim import Adam
from d3rlpy.optimizers import OptimizerFactory
from d3rlpy.algos import DQN

factory = OptimizerFactory(Adam, eps=0.001)

dqn = DQN(optim_factory=factory)
```

Parameters

- **optim_cls** – An optimizer class.

- **kwargs** – arbitrary keyword-arguments.

Methods

create (*params*, *lr*)

Returns an optimizer object.

Parameters

- **params** (*list*) – a list of PyTorch parameters.
- **lr** (*float*) – learning rate.

Returns an optimizer object.

Return type torch.optim.Optimizer

get_params (*deep=False*)

Returns optimizer parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns optimizer parameters.

Return type Dict[str, Any]

3.6.2 d3rlpy.models.optimizers.SGDFactory

class d3rlpy.models.optimizers.**SGDFactory** (*momentum=0, dampening=0, weight_decay=0, nesterov=False, **kwargs*)

An alias for SGD optimizer.

```
from d3rlpy.optimizers import SGDFactory

factory = SGDFactory(weight_decay=1e-4)
```

Parameters

- **momentum** – momentum factor.
- **dampening** – dampening for momentum.
- **weight_decay** – weight decay (L2 penalty).
- **nesterov** – flag to enable Nesterov momentum.

Methods

create (*params*, *lr*)

Returns an optimizer object.

Parameters

- **params** (*list*) – a list of PyTorch parameters.
- **lr** (*float*) – learning rate.

Returns an optimizer object.

Return type torch.optim.Optimizer

get_params (*deep=False*)

Returns optimizer parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns optimizer parameters.

Return type Dict[str, Any]

3.6.3 d3rlpy.models.optimizers.AdamFactory

class d3rlpy.models.optimizers.**AdamFactory** (*betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False, **kwargs*)

An alias for Adam optimizer.

```
from d3rlpy.optimizers import AdamFactory

factory = AdamFactory(weight_decay=1e-4)
```

Parameters

- **betas** – coefficients used for computing running averages of gradient and its square.
- **eps** – term added to the denominator to improve numerical stability.
- **weight_decay** – weight decay (L2 penalty).
- **amsgrad** – flag to use the AMSGrad variant of this algorithm.

Methods

create (*params, lr*)

Returns an optimizer object.

Parameters

- **params** (*list*) – a list of PyTorch parameters.
- **lr** (*float*) – learning rate.

Returns an optimizer object.

Return type torch.optim.Optimizer

get_params (*deep=False*)

Returns optimizer parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns optimizer parameters.

Return type Dict[str, Any]

3.6.4 d3rlpy.models.optimizers.RMSpropFactory

class d3rlpy.models.optimizers.**RMSpropFactory** (*alpha=0.95, eps=0.01, weight_decay=0, momentum=0, centered=True, **kwargs*)

An alias for RMSprop optimizer.

```
from d3rlpy.optimizers import RMSpropFactory

factory = RMSpropFactory(weight_decay=1e-4)
```

Parameters

- **alpha** – smoothing constant.
- **eps** – term added to the denominator to improve numerical stability.
- **weight_decay** – weight decay (L2 penalty).
- **momentum** – momentum factor.
- **centered** – flag to compute the centered RMSProp, the gradient is normalized by an estimation of its variance.

Methods

create (*params, lr*)

Returns an optimizer object.

Parameters

- **params** (*list*) – a list of PyTorch parameters.
- **lr** (*float*) – learning rate.

Returns an optimizer object.

Return type torch.optim.Optimizer

get_params (*deep=False*)

Returns optimizer parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns optimizer parameters.

Return type Dict[str, Any]

3.7 Network Architectures

In d3rlpy, the neural network architecture is automatically selected based on observation shape. If the observation is image, the algorithm uses the Nature DQN-based encoder at each function. Otherwise, the standard MLP architecture that consists with two linear layers with 256 hidden units.

Furthermore, d3rlpy provides `EncoderFactory` that gives you flexible control over this neural network architectures.


```

from d3rlpy.algos import DQN
from d3rlpy.models.encoders import VectorEncoderFactory

# encoder factory
encoder_factory = VectorEncoderFactory(hidden_units=[300, 400], activation='tanh')

# set EncoderFactory
dqn = DQN(encoder_factory=encoder_factory)

```

You can also build your own encoder factory.

```

import torch
import torch.nn as nn

from d3rlpy.models.encoders import EncoderFactory

# your own neural network
class CustomEncoder(nn.Module):
    def __init__(self, observation_shape, feature_size):
        self.feature_size = feature_size
        self.fc1 = nn.Linear(observation_shape[0], 64)
        self.fc2 = nn.Linear(64, feature_size)

    def forward(self, x):
        h = torch.relu(self.fc1(x))
        h = torch.relu(self.fc2(h))
        return h

    # THIS IS IMPORTANT!
    def get_feature_size(self):
        return self.feature_size

# your own encoder factory
class CustomEncoderFactory(EncoderFactory):
    TYPE = 'custom' # this is necessary

    def __init__(self, feature_size):
        self.feature_size = feature_size

    def create(self, observation_shape):
        return CustomEncoder(observation_shape, self.feature_size)

    def get_params(self, deep=False):
        return {'feature_size': self.feature_size}

dqn = DQN(encoder_factory=CustomEncoderFactory(feature_size=64))

```

You can also define action-conditioned networks such as Q-functions for continuous controls. `create` or `create_with_action` will be called depending on the function.

```

class CustomEncoderWithAction(nn.Module):
    def __init__(self, observation_shape, action_size, feature_size):
        self.feature_size = feature_size
        self.fc1 = nn.Linear(observation_shape[0] + action_size, 64)
        self.fc2 = nn.Linear(64, feature_size)

    def forward(self, x, action): # action is also given

```

(continues on next page)

(continued from previous page)

```

        h = torch.cat([x, action], dim=1)
        h = torch.relu(self.fc1(h))
        h = torch.relu(self.fc2(h))
        return h

    def get_feature_size(self):
        return self.feature_size

class CustomEncoderFactory(EncoderFactory):
    TYPE = 'custom' # this is necessary

    def __init__(self, feature_size):
        self.feature_size = feature_size

    def create(self, observation_shape):
        return CustomEncoder(observation_shape, self.feature_size)

    def create_with_action(observation_shape, action_size, discrete_action):
        return CustomEncoderWithAction(observation_shape, action_size, self.feature_
↪size)

    def get_params(self, deep=False):
        return {'feature_size': self.feature_size}

from d3rlpy.algos import SAC

factory = CustomEncoderFactory(feature_size=64)

sac = SAC(actor_encoder_factory=factory, critic_encoder_factory=factory)

```

If you want `from_json` method to load the algorithm configuration including your encoder configuration, you need to register your encoder factory.

```

from d3rlpy.models.encoders import register_encoder_factory

# register your own encoder factory
register_encoder_factory(CustomEncoderFactory)

# load algorithm from json
dqn = DQN.from_json('<path-to-json>/params.json')

```

Once you register your encoder factory, you can specify it via `TYPE` value.

```
dqn = DQN(encoder_factory='custom')
```

<code>d3rlpy.models.encoders.DefaultEncoderFactory</code>	Default encoder factory class.
<code>d3rlpy.models.encoders.PixelEncoderFactory</code>	Pixel encoder factory class.
<code>d3rlpy.models.encoders.VectorEncoderFactory</code>	Vector encoder factory class.
<code>d3rlpy.models.encoders.DenseEncoderFactory</code>	DenseNet encoder factory class.

3.7.1 d3rlpy.models.encoders.DefaultEncoderFactory

```
class d3rlpy.models.encoders.DefaultEncoderFactory (activation='relu',
                                                    use_batch_norm=False,
                                                    dropout_rate=None)
```

Default encoder factory class.

This encoder factory returns an encoder based on observation shape.

Parameters

- **activation** (*str*) – activation function name.
- **use_batch_norm** (*bool*) – flag to insert batch normalization layers.
- **dropout_rate** (*float*) – dropout probability.

Methods

create (*observation_shape*)

Returns PyTorch's state encoder module.

Parameters **observation_shape** (*Sequence[int]*) – observation shape.

Returns an encoder object.

Return type d3rlpy.models.torch.encoders.Encoder

create_with_action (*observation_shape, action_size, discrete_action=False*)

Returns PyTorch's state-action encoder module.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – action size. If None, the encoder does not take action as input.
- **discrete_action** (*bool*) – flag if action-space is discrete.

Returns an encoder object.

Return type d3rlpy.models.torch.encoders.EncoderWithAction

get_params (*deep=False*)

Returns encoder parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns encoder parameters.

Return type Dict[str, Any]

get_type ()

Returns encoder type.

Returns encoder type.

Return type str

Attributes

TYPE: `ClassVar[str] = 'default'`

3.7.2 d3rlpy.models.encoders.PixelEncoderFactory

```
class d3rlpy.models.encoders.PixelEncoderFactory (filters=None,           fea-  
                                                    ture_size=512,    activation='relu',  
                                                    use_batch_norm=False,  
                                                    dropout_rate=None)
```

Pixel encoder factory class.

This is the default encoder factory for image observation.

Parameters

- **filters** (*list*) – list of tuples consisting with (*filter_size*, *kernel_size*, *stride*). If None, Nature DQN-based architecture is used.
- **feature_size** (*int*) – the last linear layer size.
- **activation** (*str*) – activation function name.
- **use_batch_norm** (*bool*) – flag to insert batch normalization layers.
- **dropout_rate** (*float*) – dropout probability.

Methods

create (*observation_shape*)

Returns PyTorch's state encoder module.

Parameters **observation_shape** (*Sequence[int]*) – observation shape.

Returns an encoder object.

Return type `d3rlpy.models.torch.encoders.PixelEncoder`

create_with_action (*observation_shape*, *action_size*, *discrete_action=False*)

Returns PyTorch's state-action encoder module.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – action size. If None, the encoder does not take action as input.
- **discrete_action** (*bool*) – flag if action-space is discrete.

Returns an encoder object.

Return type `d3rlpy.models.torch.encoders.PixelEncoderWithAction`

get_params (*deep=False*)

Returns encoder parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns encoder parameters.

Return type `Dict[str, Any]`

get_type ()

Returns encoder type.

Returns encoder type.

Return type `str`

Attributes

TYPE: `ClassVar[str] = 'pixel'`

3.7.3 d3rlpy.models.encoders.VectorEncoderFactory

```
class d3rlpy.models.encoders.VectorEncoderFactory(hidden_units=None, activation='relu', use_batch_norm=False, dropout_rate=None, use_dense=False)
```

Vector encoder factory class.

This is the default encoder factory for vector observation.

Parameters

- **hidden_units** (*list*) – list of hidden unit sizes. If `None`, the standard architecture with `[256, 256]` is used.
- **activation** (*str*) – activation function name.
- **use_batch_norm** (*bool*) – flag to insert batch normalization layers.
- **use_dense** (*bool*) – flag to use DenseNet architecture.
- **dropout_rate** (*float*) – dropout probability.

Methods

create (*observation_shape*)

Returns PyTorch's state encoder module.

Parameters **observation_shape** (*Sequence[int]*) – observation shape.

Returns an encoder object.

Return type `d3rlpy.models.torch.encoders.VectorEncoder`

create_with_action (*observation_shape, action_size, discrete_action=False*)

Returns PyTorch's state-action encoder module.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – action size. If `None`, the encoder does not take action as input.
- **discrete_action** (*bool*) – flag if action-space is discrete.

Returns an encoder object.

Return type `d3rlpy.models.torch.encoders.VectorEncoderWithAction`

get_params (*deep=False*)

Returns encoder parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns encoder parameters.

Return type Dict[str, Any]

get_type()

Returns encoder type.

Returns encoder type.

Return type str

Attributes

TYPE: ClassVar[str] = 'vector'

3.7.4 d3rlpy.models.encoders.DenseEncoderFactory

```
class d3rlpy.models.encoders.DenseEncoderFactory(activation='relu',  
                                                  use_batch_norm=False,  
                                                  dropout_rate=None)
```

DenseNet encoder factory class.

This is an alias for DenseNet architecture proposed in D2RL. This class does exactly same as follows.

```
from d3rlpy.encoders import VectorEncoderFactory  
  
factory = VectorEncoderFactory(hidden_units=[256, 256, 256, 256],  
                               use_dense=True)
```

For now, this only supports vector observations.

References

- Sinha et al., D2RL: Deep Dense Architectures in Reinforcement Learning.

Parameters

- **activation** (str) – activation function name.
- **use_batch_norm** (bool) – flag to insert batch normalization layers.
- **dropout_rate** (float) – dropout probability.

Methods

create (observation_shape)

Returns PyTorch's state encoder module.

Parameters **observation_shape** (Sequence[int]) – observation shape.

Returns an encoder object.

Return type d3rlpy.models.torch.encoders.VectorEncoder

create_with_action (observation_shape, action_size, discrete_action=False)

Returns PyTorch's state-action encoder module.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – action size. If None, the encoder does not take action as input.
- **discrete_action** (*bool*) – flag if action-space is discrete.

Returns an encoder object.

Return type `d3rlpy.models.torch.encoders.VectorEncoderWithAction`

get_params (*deep=False*)

Returns encoder parameters.

Parameters **deep** (*bool*) – flag to deeply copy the parameters.

Returns encoder parameters.

Return type `Dict[str, Any]`

get_type ()

Returns encoder type.

Returns encoder type.

Return type `str`

Attributes

TYPE: ClassVar[`str`] = 'dense'

3.8 Metrics

d3rlpy provides scoring functions without compromising scikit-learn compatibility. You can evaluate many metrics with test episodes during training.

```
from d3rlpy.datasets import get_cartpole
from d3rlpy.algos import DQN
from d3rlpy.metrics.scorer import td_error_scorer
from d3rlpy.metrics.scorer import average_value_estimation_scorer
from d3rlpy.metrics.scorer import evaluate_on_environment
from sklearn.model_selection import train_test_split

dataset, env = get_cartpole()

train_episodes, test_episodes = train_test_split(dataset)

dqn = DQN()

dqn.fit(train_episodes,
        eval_episodes=test_episodes,
        scorers={
            'td_error': td_error_scorer,
            'value_scale': average_value_estimation_scorer,
            'environment': evaluate_on_environment(env)
        })
```

You can also use them with scikit-learn utilities.

```

from sklearn.model_selection import cross_validate

scores = cross_validate(dqn,
                        dataset,
                        scoring={
                            'td_error': td_error_scorer,
                            'environment': evaluate_on_environment(env)
                        })

```

3.8.1 Algorithms

<code>d3rlpy.metrics.scorer.td_error_scorer</code>	Returns average TD error (in negative scale).
<code>d3rlpy.metrics.scorer.discounted_sum_of_advantage_scorer</code>	Returns average of discounted sum of advantage (in negative scale).
<code>d3rlpy.metrics.scorer.average_value_estimation_scorer</code>	Returns average value estimation (in negative scale).
<code>d3rlpy.metrics.scorer.value_estimation_std_scorer</code>	Returns standard deviation of value estimation (in negative scale).
<code>d3rlpy.metrics.scorer.initial_state_value_estimation_scorer</code>	Returns mean estimated action-values at the initial states.
<code>d3rlpy.metrics.scorer.soft_opc_scorer</code>	Returns Soft Off-Policy Classification metrics.
<code>d3rlpy.metrics.scorer.continuous_action_diff_scorer</code>	Returns squared difference of actions between algorithm and dataset.
<code>d3rlpy.metrics.scorer.discrete_action_match_scorer</code>	Returns percentage of identical actions between algorithm and dataset.
<code>d3rlpy.metrics.scorer.evaluate_on_environment</code>	Returns scorer function of evaluation on environment.
<code>d3rlpy.metrics.comparer.compare_continuous_action_diff</code>	Returns scorer function of action difference between algorithms.
<code>d3rlpy.metrics.comparer.compare_discrete_action_match</code>	Returns scorer function of action matches between algorithms.

d3rlpy.metrics.scorer.td_error_scorer

`d3rlpy.metrics.scorer.td_error_scorer` (*algo, episodes*)

Returns average TD error (in negative scale).

This metrics suggests how Q functions overfit to training sets. If the TD error is large, the Q functions are overfitting.

$$\mathbb{E}_{s_t, a_t, r_{t+1}, s_{t+1} \sim D} [Q_{\theta}(s_t, a_t) - (r_t + \gamma \max_a Q_{\theta}(s_{t+1}, a))^2]$$

Parameters

- **algo** (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative average TD error.

Return type `float`

d3rlpy.metrics.scorer.discounted_sum_of_advantage_scorer

`d3rlpy.metrics.scorer.discounted_sum_of_advantage_scorer` (*algo*, *episodes*)

Returns average of discounted sum of advantage (in negative scale).

This metrics suggests how the greedy-policy selects different actions in action-value space. If the sum of advantage is small, the policy selects actions with larger estimated action-values.

$$\mathbb{E}_{s_t, a_t \sim D} \left[\sum_{t'=t} \gamma^{t'-t} A(s_{t'}, a_{t'}) \right]$$

where $A(s_t, a_t) = Q_\theta(s_t, a_t) - \max_a Q_\theta(s_t, a)$.

References

- [Murphy, A generalization error for Q-Learning.](#)

Parameters

- **algo** (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative average of discounted sum of advantage.

Return type `float`

d3rlpy.metrics.scorer.average_value_estimation_scorer

`d3rlpy.metrics.scorer.average_value_estimation_scorer` (*algo*, *episodes*)

Returns average value estimation (in negative scale).

This metrics suggests the scale for estimation of Q functions. If average value estimation is too large, the Q functions overestimate action-values, which possibly makes training failed.

$$\mathbb{E}_{s_t \sim D} [\max_a Q_\theta(s_t, a)]$$

Parameters

- **algo** (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative average value estimation.

Return type `float`

d3rlpy.metrics.scorer.value_estimation_std_scorer

`d3rlpy.metrics.scorer.value_estimation_std_scorer` (*algo*, *episodes*)

Returns standard deviation of value estimation (in negative scale).

This metrics suggests how confident Q functions are for the given episodes. This metrics will be more accurate with *bootstrap* enabled and the larger *n_critics* at algorithm. If standard deviation of value estimation is large, the Q functions are overfitting to the training set.

$$\mathbb{E}_{s_t \sim D, a \sim \arg\max_a Q_\theta(s_t, a)} [Q_{\text{std}}(s_t, a)]$$

where $Q_{\text{std}}(s, a)$ is a standard deviation of action-value estimation over ensemble functions.

Parameters

- **algo** (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative standard deviation.

Return type `float`

d3rlpy.metrics.scorer.initial_state_value_estimation_scorer

`d3rlpy.metrics.scorer.initial_state_value_estimation_scorer` (*algo, episodes*)

Returns mean estimated action-values at the initial states.

This metrics suggests how much return the trained policy would get from the initial states by deploying the policy to the states. If the estimated value is large, the trained policy is expected to get higher returns.

$$\mathbb{E}_{s_0 \sim D}[Q(s_0, \pi(s_0))]$$

References

- Paine et al., Hyperparameter Selection for Offline Reinforcement Learning

Parameters

- **algo** (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns mean action-value estimation at the initial states.

Return type `float`

d3rlpy.metrics.scorer.soft_opc_scorer

`d3rlpy.metrics.scorer.soft_opc_scorer` (*return_threshold*)

Returns Soft Off-Policy Classification metrics.

This function returns scorer function, which is suitable to the standard scikit-learn scorer function style. The metrics of the scorer function is evaluating gaps of action-value estimation between the success episodes and the all episodes. If the learned Q-function is optimal, action-values in success episodes are expected to be higher than the others. The success episode is defined as an episode with a return above the given threshold.

$$\mathbb{E}_{s,a \sim D_{success}}[Q(s,a)] - \mathbb{E}_{s,a \sim D}[Q(s,a)]$$

```
from d3rlpy.datasets import get_cartpole
from d3rlpy.algos import DQN
from d3rlpy.metrics.scorer import soft_opc_scorer
from sklearn.model_selection import train_test_split

dataset, _ = get_cartpole()
train_episodes, test_episodes = train_test_split(dataset, test_size=0.2)

scorer = soft_opc_scorer(return_threshold=180)
```

(continues on next page)

(continued from previous page)

```
dqn = DQN()
dqn.fit(train_episodes,
        eval_episodes=test_episodes,
        scorers={'soft_opc': scorer})
```

References

- Irpan et al., Off-Policy Evaluation via Off-Policy Classification.

Parameters `return_threshold` (*float*) – threshold of success episodes.

Returns scorer function.

Return type Callable[[d3rlpy.metrics.scorer.AlgoProtocol, List[d3rlpy.dataset.Episode]], float]

d3rlpy.metrics.scorer.continuous_action_diff_scorer

`d3rlpy.metrics.scorer.continuous_action_diff_scorer` (*algo*, *episodes*)

Returns squared difference of actions between algorithm and dataset.

This metrics suggests how different the greedy-policy is from the given episodes in continuous action-space. If the given episodes are near-optimal, the small action difference would be better.

$$\mathbb{E}_{s_t, a_t \sim D}[(a_t - \pi_\phi(s_t))^2]$$

Parameters

- **algo** (*d3rlpy.metrics.scorer.AlgoProtocol*) – algorithm.
- **episodes** (*List[d3rlpy.dataset.Episode]*) – list of episodes.

Returns negative squared action difference.

Return type float

d3rlpy.metrics.scorer.discrete_action_match_scorer

`d3rlpy.metrics.scorer.discrete_action_match_scorer` (*algo*, *episodes*)

Returns percentage of identical actions between algorithm and dataset.

This metrics suggests how different the greedy-policy is from the given episodes in discrete action-space. If the given episodes are near-optimal, the large percentage would be better.

$$\frac{1}{N} \sum^N \mathbb{I}\{a_t = \operatorname{argmax}_a Q_\theta(s_t, a)\}$$

Parameters

- **algo** (*d3rlpy.metrics.scorer.AlgoProtocol*) – algorithm.
- **episodes** (*List[d3rlpy.dataset.Episode]*) – list of episodes.

Returns percentage of identical actions.

Return type float

d3rlpy.metrics.scorer.evaluate_on_environment

`d3rlpy.metrics.scorer.evaluate_on_environment` (*env*, *n_trials=10*, *epsilon=0.0*, *render=False*)

Returns scorer function of evaluation on environment.

This function returns scorer function, which is suitable to the standard scikit-learn scorer function style. The metrics of the scorer function is ideal metrics to evaluate the resulted policies.

```
import gym

from d3rlpy.algos import DQN
from d3rlpy.metrics.scorer import evaluate_on_environment

env = gym.make('CartPole-v0')

scorer = evaluate_on_environment(env)

cql = CQL()

mean_episode_return = scorer(cql)
```

Parameters

- **env** (*gym.core.Env*) – gym-styled environment.
- **n_trials** (*int*) – the number of trials.
- **epsilon** (*float*) – noise factor for epsilon-greedy policy.
- **render** (*bool*) – flag to render environment.

Returns scorer function.

Return type Callable[[*...*], *float*]

d3rlpy.metrics.comparer.compare_continuous_action_diff

`d3rlpy.metrics.comparer.compare_continuous_action_diff` (*base_algo*)

Returns scorer function of action difference between algorithms.

This metrics suggests how different the two algorithms are in continuous action-space. If the algorithm to compare with is near-optimal, the small action difference would be better.

$$\mathbb{E}_{s_t \sim D}[(\pi_{\phi_1}(s_t) - \pi_{\phi_2}(s_t))^2]$$

```
from d3rlpy.algos import CQL
from d3rlpy.metrics.comparer import compare_continuous_action_diff

cql1 = CQL()
cql2 = CQL()

scorer = compare_continuous_action_diff(cql1)

squared_action_diff = scorer(cql2, ...)
```

Parameters `base_algo` (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm to compare with.

Returns scorer function.

Return type Callable[[`d3rlpy.metrics.scorer.AlgoProtocol`, List[`d3rlpy.dataset.Episode`]], float]

d3rlpy.metrics.comparer.compare_discrete_action_match

`d3rlpy.metrics.comparer.compare_discrete_action_match` (`base_algo`)

Returns scorer function of action matches between algorithms.

This metrics suggests how different the two algorithms are in discrete action-space. If the algorithm to compare with is near-optimal, the small action difference would be better.

$$\mathbb{E}_{s_t \sim D} [\| \{ \operatorname{argmax}_a Q_{\theta_1}(s_t, a) = \operatorname{argmax}_a Q_{\theta_2}(s_t, a) \} \|]$$

```
from d3rlpy.algos import DQN
from d3rlpy.metrics.comparer import compare_continuous_action_diff

dqn1 = DQN()
dqn2 = DQN()

scorer = compare_continuous_action_diff(dqn1)

percentage_of_identical_actions = scorer(dqn2, ...)
```

Parameters `base_algo` (`d3rlpy.metrics.scorer.AlgoProtocol`) – algorithm to compare with.

Returns scorer function.

Return type Callable[[`d3rlpy.metrics.scorer.AlgoProtocol`, List[`d3rlpy.dataset.Episode`]], float]

3.8.2 Dynamics

<code>d3rlpy.metrics.scorer.dynamics_observation_prediction_error_scorer</code>	Returns MSE of observation prediction (in negative scale).
<code>d3rlpy.metrics.scorer.dynamics_reward_prediction_error_scorer</code>	Returns MSE of reward prediction (in negative scale).
<code>d3rlpy.metrics.scorer.dynamics_prediction_variance_scorer</code>	Returns prediction variance of ensemble dynamics (in negative scale).

d3rlpy.metrics.scorer.dynamics_observation_prediction_error_scorer

`d3rlpy.metrics.scorer.dynamics_observation_prediction_error_scorer` (*dynamics*,
episodes)

Returns MSE of observation prediction (in negative scale).

This metrics suggests how dynamics model is generalized to test sets. If the MSE is large, the dynamics model are overfitting.

$$\mathbb{E}_{s_t, a_t, s_{t+1} \sim D} [(s_{t+1} - s')^2]$$

where $s' \sim T(s_t, a_t)$.

Parameters

- **dynamics** (`d3rlpy.metrics.scorer.DynamicsProtocol`) – dynamics model.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative mean squared error.

Return type `float`

d3rlpy.metrics.scorer.dynamics_reward_prediction_error_scorer

`d3rlpy.metrics.scorer.dynamics_reward_prediction_error_scorer` (*dynamics*,
episodes)

Returns MSE of reward prediction (in negative scale).

This metrics suggests how dynamics model is generalized to test sets. If the MSE is large, the dynamics model are overfitting.

$$\mathbb{E}_{s_t, a_t, r_{t+1} \sim D} [(r_{t+1} - r')^2]$$

where $r' \sim T(s_t, a_t)$.

Parameters

- **dynamics** (`d3rlpy.metrics.scorer.DynamicsProtocol`) – dynamics model.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative mean squared error.

Return type `float`

d3rlpy.metrics.scorer.dynamics_prediction_variance_scorer

`d3rlpy.metrics.scorer.dynamics_prediction_variance_scorer` (*dynamics*, *episodes*)

Returns prediction variance of ensemble dynamics (in negative scale).

This metrics suggests how dynamics model is confident of test sets. If the variance is large, the dynamics model has large uncertainty.

Parameters

- **dynamics** (`d3rlpy.metrics.scorer.DynamicsProtocol`) – dynamics model.
- **episodes** (`List[d3rlpy.dataset.Episode]`) – list of episodes.

Returns negative variance.

Return type `float`

3.9 Off-Policy Evaluation

The off-policy evaluation is a method to estimate the trained policy performance only with offline datasets.

```
from d3rlpy.algos import CQL
from d3rlpy.datasets import get_pybullet

# prepare the trained algorithm
cql = CQL.from_json('<path-to-json>/params.json')
cql.load_model('<path-to-model>/model.pt')

# dataset to evaluate with
dataset, env = get_pybullet('hopper-bullet-mixed-v0')

from d3rlpy.ope import FQE

# off-policy evaluation algorithm
fqe = FQE(algo=cql)

# metrics to evaluate with
from d3rlpy.metrics.scorer import initial_state_value_estimation_scorer
from d3rlpy.metrics.scorer import soft_opc_scorer

# train estimators to evaluate the trained policy
fqe.fit(dataset.episodes,
        eval_episodes=dataset.episodes,
        scorers={
            'init_value': initial_state_value_estimation_scorer,
            'soft_opc': soft_opc_scorer(return_threshold=600)
        })
```

The evaluation during fitting is evaluating the trained policy.

3.9.1 For continuous control algorithms

`d3rlpy.ope.FQE`

Fitted Q Evaluation.

`d3rlpy.ope.FQE`

```
class d3rlpy.ope.FQE(*,
                    algo=None,
                    learning_rate=0.0001,
                    optim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                betas=(0.9, 0.999), eps=1e-08, weight_decay=0, amsgrad=False),
                    encoder_factory='default', q_func_factory='mean', batch_size=100, n_frames=1,
                    n_steps=1, gamma=0.99, n_critics=1, target_update_interval=100,
                    use_gpu=False, scaler=None, action_scaler=None, impl=None, **kwargs)
```

Fitted Q Evaluation.

FQE is an off-policy evaluation method that approximates a Q function $Q_{\theta}(s, a)$ with the trained policy $\pi_{\phi}(s)$.

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1} | s_{t+1} \sim D} [(Q_{\theta}(s_t, a_t) - r_{t+1} - \gamma Q_{\theta'}(s_{t+1}, \pi_{\phi}(s_{t+1})))^2]$$

The trained Q function in FQE will estimate evaluation metrics more accurately than learned Q function during training.

References

- Le et al., Batch Policy Learning under Constraints.

Parameters

- **algo** (`d3rlpy.algos.base.AlgoBase`) – algorithm to evaluate.
- **learning_rate** (`float`) – learning rate.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory` or `str`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_update_interval** (`int`) – interval to update the target network.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or `str`) – preprocessor. The available options are [`'pixel'`, `'min_max'`, `'standard'`].
- **action_scaler** (`d3rlpy.preprocessing.ActionScaler` or `str`) – action preprocessor. The available options are [`'min_max'`].
- **impl** (`d3rlpy.metrics.ope.torch.FQEImpl`) – algorithm implementation.

Methods

build_with_dataset (`dataset`)

Instantiate implementation object with MDPDataset object.

Parameters `dataset` (`d3rlpy.dataset.MDPDataset`) – dataset.

Return type `None`

build_with_env (`env`)

Instantiate implementation object with OpenAI Gym object.

Parameters `env` (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (`env`, `buffer=None`, `explorer=None`, `n_steps=1000000`, `show_progress=True`, `time_limit_aware=True`)

Collects data via interaction with environment.

If `buffer` is not given, `ReplayBuffer` will be internally created.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape, action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.

- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fit_batch_online (*env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000, n_updates_per_epoch=1000, eval_interval=10, eval_env=None, eval_epsilon=0.0, save_metrics=True, save_interval=1, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, timelimit_aware=True, callback=None*)

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.

- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Return type *None*

fit_online (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *n_steps_per_epoch=10000*, *update_interval=1*, *update_start_step=0*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.

- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (`algo`, `epoch`, `total_step`), which is called at the end of epochs.

Return type `None`

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when `n_steps` is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with `eval_episodes`.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.

- **callback** (*Optional*[*Callable*[[*d3rlpy.base.LearnableBase*, *int*, *int*], *None*]]) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type *Generator*[*Tuple*[*int*, *Dict*[*str*, *float*]], *None*, *None*]

classmethod from_json (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional*[*Union*[*bool*, *int*, *d3rlpy.gpu.Device*]]) – flag to use GPU, device ID or device.

Returns algorithm.

Return type *d3rlpy.base.LearnableBase*

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List*[*d3rlpy.dataset.Transition*]) – list of transitions.

Returns list of new transitions.

Return type *Optional*[*List*[*d3rlpy.dataset.Transition*]]

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type *d3rlpy.constants.ActionSpace*

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters `deep` (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters `fname` (*str*) – source file path.

Return type None

predict (*x*)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters `x` (*Union[numpy.ndarray, List[Any]]*) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (*x, action, with_std=False*)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (`Union[numpy.ndarray, List[Any]]`) – observations
- **action** (`Union[numpy.ndarray, List[Any]]`) – actions
- **with_std** (`bool`) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type `Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]`

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters **x** (`Union[numpy.ndarray, List[Any]]`) – observations.

Returns sampled actions.

Return type `numpy.ndarray`

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type `None`

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (`d3rlpy.logger.D3RLPyLogger`) – logger object.

Return type `None`

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- **fname** (*str*) – destination file path.
- **as_onnx** (*bool*) – flag to save as ONNX format.

Return type `None`

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type `d3rlpy.base.LearnableBase`

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type `Dict[str, float]`

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type `Optional[ActionScaler]`

action_size

Action size.

Returns action size.

Return type `Optional[int]`

batch_size

Batch size to train.

Returns batch size.

Return type `int`

gamma

Discount factor.

Returns discount factor.

Return type `float`

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

3.9.2 For discrete control algorithms

`d3rlpy.ope.DiscreteFQE`

Fitted Q Evaluation for discrete action-space.

`d3rlpy.ope.DiscreteFQE`

```
class d3rlpy.ope.DiscreteFQE(*,          algo=None,          learning_rate=0.0001,          op-
                                tim_factory=d3rlpy.models.optimizers.AdamFactory(optim_cls='Adam',
                                betas=(0.9, 0.999), eps=1e-08, weight_decay=0, ams-
                                grad=False), encoder_factory='default', q_func_factory='mean',
                                batch_size=100, n_frames=1, n_steps=1, gamma=0.99,
                                n_critics=1, target_update_interval=100, use_gpu=False,
                                scaler=None, action_scaler=None, impl=None, **kwargs)
```

Fitted Q Evaluation for discrete action-space.

FQE is an off-policy evaluation method that approximates a Q function $Q_\theta(s, a)$ with the trained policy $\pi_\phi(s)$.

$$L(\theta) = \mathbb{E}_{s_t, a_t, r_{t+1} | s_{t+1} \sim D} [(Q_\theta(s_t, a_t) - r_{t+1} - \gamma Q_{\theta'}(s_{t+1}, \pi_\phi(s_{t+1})))^2]$$

The trained Q function in FQE will estimate evaluation metrics more accurately than learned Q function during training.

References

- Le et al., Batch Policy Learning under Constraints.

Parameters

- **algo** (`d3rlpy.algos.base.AlgoBase`) – algorithm to evaluate.
- **learning_rate** (`float`) – learning rate.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory` or `str`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or `str`) – encoder factory.
- **q_func_factory** (`d3rlpy.models.q_functions.QFunctionFactory` or `str`) – Q function factory.
- **batch_size** (`int`) – mini-batch size.
- **n_frames** (`int`) – the number of frames to stack for image observation.
- **n_steps** (`int`) – N-step TD calculation.
- **gamma** (`float`) – discount factor.
- **n_critics** (`int`) – the number of Q functions for ensemble.
- **target_update_interval** (`int`) – interval to update the target network.
- **use_gpu** (`bool`, `int` or `d3rlpy.gpu.Device`) – flag to use GPU, device ID or device.
- **scaler** (`d3rlpy.preprocessingScaler` or `str`) – preprocessor. The available options are [`'pixel'`, `'min_max'`, `'standard'`]
- **impl** (`d3rlpy.metrics.opentorch.FQEImpl`) – algorithm implementation.

Methods

build_with_dataset (`dataset`)

Instantiate implementation object with `MDPDataSet` object.

Parameters `dataset` (`d3rlpy.dataset.MDPDataSet`) – dataset.

Return type `None`

build_with_env (`env`)

Instantiate implementation object with OpenAI Gym object.

Parameters `env` (`gym.core.Env`) – gym-like environment.

Return type `None`

collect (`env`, `buffer=None`, `explorer=None`, `n_steps=1000000`, `show_progress=True`, `time-limit_aware=True`)

Collects data via interaction with environment.

If `buffer` is not given, `ReplayBuffer` will be internally created.

Parameters

- `env` (`gym.core.Env`) – gym-like environment.

- **buffer** (*Optional[d3rlpy.online.buffers.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **timelimit_aware** (*bool*) – flag to turn terminal flag `False` when `TimeLimit.truncated` flag is `True`, which is designed to incorporate with `gym.wrappers.TimeLimit`.

Returns replay buffer with the collected data.

Return type `d3rlpy.online.buffers.Buffer`

create_impl (*observation_shape, action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None*)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is `None`.
- **save_metrics** (*bool*) – flag to record metrics in files. If `False`, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be `{class name}_{timestamp}`.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if `None`, the directory will not be created.

- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

```
fit_batch_online (env, buffer=None, explorer=None, n_epochs=1000, n_steps_per_epoch=1000,  
                  n_updates_per_epoch=1000, eval_interval=10, eval_env=None,  
                  eval_epsilon=0.0, save_metrics=True, save_interval=1, experi-  
                  ment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True,  
                  show_progress=True, tensorboard_dir=None, timelimit_aware=True, call-  
                  back=None)
```

Start training loop of batch online deep reinforcement learning.

Parameters

- **env** (*d3rlpy.envs.batch.BatchEnv*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffers.BatchBuffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_epochs** (*int*) – the number of epochs to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** – the number of steps per update.
- **n_updates_per_epoch** (*int*) – the number of updates per epoch.
- **eval_interval** (*int*) – the number of epochs before evaluation.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If None, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If False, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.

- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.
- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type *None*

fit_online (*env*, *buffer=None*, *explorer=None*, *n_steps=1000000*, *n_steps_per_epoch=10000*, *update_interval=1*, *update_start_step=0*, *eval_env=None*, *eval_epsilon=0.0*, *save_metrics=True*, *save_interval=1*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *timelimit_aware=True*, *callback=None*)

Start training loop of online deep reinforcement learning.

Parameters

- **env** (*gym.core.Env*) – gym-like environment.
- **buffer** (*Optional[d3rlpy.online.buffer.Buffer]*) – replay buffer.
- **explorer** (*Optional[d3rlpy.online.explorers.Explorer]*) – action explorer.
- **n_steps** (*int*) – the number of total steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch.
- **update_interval** (*int*) – the number of steps per update.
- **update_start_step** (*int*) – the steps before starting updates.
- **eval_env** (*Optional[gym.core.Env]*) – gym-like environment. If *None*, evaluation is skipped.
- **eval_epsilon** (*float*) – ϵ -greedy factor during evaluation.
- **save_metrics** (*bool*) – flag to record metrics. If *False*, the log directory is not created and the model parameters are not saved.
- **save_interval** (*int*) – the number of epochs before saving models.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_online_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **timelimit_aware** (*bool*) – flag to turn terminal flag *False* when *TimeLimit.truncated* flag is *True*, which is designed to incorporate with *gym.wrappers.TimeLimit*.

- **callback** (*Optional[Callable[[d3rlpy.online.iterators.AlgoProtocol, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Return type `None`

fitter (dataset, n_epochs=None, n_steps=None, n_steps_per_epoch=10000, save_metrics=True, experiment_name=None, with_timestamp=True, logdir='d3rlpy_logs', verbose=True, show_progress=True, tensorboard_dir=None, eval_episodes=None, save_interval=1, scorers=None, shuffle=True, callback=None)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when n_steps is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be {class name}_{timestamp}.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode], float]]]]*) – list of scorer functions used with eval_episodes.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type `Generator[Tuple[int, Dict[str, float]], None, None]`

classmethod `from_json` (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (fname)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (str) – source file path.

Return type None

predict (x)

Returns greedy actions.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

actions = algo.predict(x)
# actions.shape == (100, action size) for continuous control
# actions.shape == (100,) for discrete control
```

Parameters **x** (Union[numpy.ndarray, List[Any]]) – observations

Returns greedy actions

Return type numpy.ndarray

predict_value (x, action, with_std=False)

Returns predicted action-values.

```
# 100 observations with shape of (10,)
x = np.random.random((100, 10))

# for continuous control
# 100 actions with shape of (2,)
actions = np.random.random((100, 2))

# for discrete control
# 100 actions in integer values
actions = np.random.randint(2, size=100)

values = algo.predict_value(x, actions)
# values.shape == (100,)

values, stds = algo.predict_value(x, actions, with_std=True)
# stds.shape == (100,)
```

Parameters

- **x** (Union[numpy.ndarray, List[Any]]) – observations
- **action** (Union[numpy.ndarray, List[Any]]) – actions
- **with_std** (bool) – flag to return standard deviation of ensemble estimation. This deviation reflects uncertainty for the given observations. This uncertainty will be more accurate if you enable `bootstrap` flag and increase `n_critics` value.

Returns predicted action-values

Return type Union[numpy.ndarray, Tuple[numpy.ndarray, numpy.ndarray]]

sample_action (*x*)

Returns sampled actions.

The sampled actions are identical to the output of *predict* method if the policy is deterministic.

Parameters *x* (Union[numpy.ndarray, List[Any]]) – observations.

Returns sampled actions.

Return type numpy.ndarray

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters *fname* (*str*) – destination file path.

Return type None

save_params (*logger*)

Saves configurations as params.json.

Parameters *logger* (d3rlpy.logger.D3RLPyLogger) – logger object.

Return type None

save_policy (*fname*, *as_onnx=False*)

Save the greedy-policy computational graph as TorchScript or ONNX.

```
# save as TorchScript
algo.save_policy('policy.pt')

# save as ONNX
algo.save_policy('policy.onnx', as_onnx=True)
```

The artifacts saved with this method will work without d3rlpy. This method is especially useful to deploy the learned policy to production environments or embedding systems.

See also

- https://pytorch.org/tutorials/beginner/Intro_to_TorchScript_tutorial.html (for Python).
- https://pytorch.org/tutorials/advanced/cpp_export.html (for C++).
- <https://onnx.ai> (for ONNX)

Parameters

- *fname* (*str*) – destination file path.
- *as_onnx* (*bool*) – flag to save as ONNX format.

Return type None

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type d3rlpy.base.LearnableBase

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (d3rlpy.dataset.TransitionMiniBatch) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type `int`

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type `int`

observation_shape

Observation shape.

Returns observation shape.

Return type `Optional[Sequence[int]]`

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type `Optional[Scaler]`

3.10 Save and Load

3.10.1 save_model and load_model

```
from d3rlpy.datasets import get_cartpole
from d3rlpy.algos import DQN

dataset, env = get_cartpole()

dqn = DQN()
dqn.fit(dataset.episodes, n_epochs=1)

# save entire model parameters.
dqn.save_model('model.pt')
```

`save_model` method saves all parameters including optimizer states, which is useful when checking all the outputs or re-training from snapshots.

Once you save your model, you can load it via `load_model` method. Before loading the model, the algorithm object must be initialized as follows.

```
dqn = DQN()

# initialize with dataset
dqn.build_with_dataset(dataset)

# initialize with environment
# dqn.build_with_env(env)
```

(continues on next page)

```
# load entire model parameters.
dqn.load_model('model.pt')
```

3.10.2 from_json

It is very boring to set the same hyperparameters to initialize algorithms when loading model parameters. In d3rlpy, `params.json` is saved at the beginning of `fit` method, which includes all hyperparameters within the algorithm object. You can recreate algorithm objects from `params.json` via `from_json` method.

```
from d3rlpy.algos import DQN

dqn = DQN.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
dqn.load_model('model.pt')
```

3.10.3 save_policy

`save_policy` method saves the only greedy-policy computation graph as TorchScript or ONNX. When `save_policy` method is called, the greedy-policy graph is constructed and traced via `torch.jit.trace` function.

```
from d3rlpy.datasets import get_cartpole
from d3rlpy.algos import DQN

dataset, env = get_cartpole()

dqn = DQN()
dqn.fit(dataset.episodes, n_epochs=1)

# save greedy-policy as TorchScript
dqn.save_policy('policy.pt')

# save greedy-policy as ONNX
dqn.save_policy('policy.onnx', as_onnx=True)
```

TorchScript

TorchScript is a optimizable graph expression provided by PyTorch. The saved policy can be loaded without any dependencies except PyTorch.

```
import torch

# load greedy-policy only with PyTorch
policy = torch.jit.load('policy.pt')

# returns greedy actions
actions = policy(torch.rand(32, 6))
```

This is especially useful when deploying the trained models to productions. The computation can be faster and you don't need to install d3rlpy. Moreover, TorchScript model can be easily loaded even with C++, which will empower your robotics and embedding system projects.

```
#include <torch/script.h>

int main(int argc, char* argv[]) {
    torch::jit::script::Module module;
    try {
        module = torch::jit::load("policy.pt")
    } catch (const c10::Error& e) {
        return -1;
    }
    return 0;
}
```

You can get more information about TorchScript [here](#).

ONNX

ONNX is an open format built to represent machine learning models. This is also useful when deploying the trained model to productions with various programming languages including Python, C++, JavaScript and more.

The following example is written with `onnxruntime`.

```
import onnxruntime as ort

# load ONNX policy via onnxruntime
ort_session = ort.InferenceSession('policy.onnx')

# observation
observation = np.random.rand(1, 6).astype(np.float32)

# returns greedy action
action = ort_session.run(None, {'input_0': observation})[0]
```

You can get more information about ONNX [here](#).

3.11 Logging

d3rlpy algorithms automatically save model parameters and metrics under `d3rlpy_logs` directory.

```
from d3rlpy.datasets import get_cartpole
from d3rlpy.algos import DQN

dataset, env = get_cartpole()

dqn = DQN()

# metrics and parameters are saved in `d3rlpy_logs/DQN_YYYYMMDDHHmmss`
dqn.fit(dataset.episodes)
```

You can designate the directory.

```
# the directory will be `custom_logs/custom_YYYYMMDDHHmmss`
dqn.fit(dataset.episodes, logdir='custom_logs', experiment_name='custom')
```

If you want to disable all loggings, you can pass `save_metrics=False`.

```
dqn.fit(dataset.episodes, save_metrics=False)
```

3.11.1 TensorBoard

The same information is also automatically saved for tensorboard under *runs* directory. You can interactively visualize training metrics easily.

```
$ pip install tensorboard
$ tensorboard --logdir runs
```

This tensorboard logs can be enabled by passing *tensorboard_dir=/path/to/log_dir*.

```
dqn.fit(dataset.episodes, tensorboard_dir='runs')
```

3.12 scikit-learn compatibility

d3rlpy provides complete scikit-learn compatible APIs.

3.12.1 train_test_split

d3rlpy.dataset.MDPDataset is compatible with splitting functions in scikit-learn.

```
from d3rlpy.algos import DQN
from d3rlpy.datasets import get_cartpole
from d3rlpy.metrics.scorer import td_error_scorer
from sklearn.model_selection import train_test_split

dataset, env = get_cartpole()

train_episodes, test_episodes = train_test_split(dataset, test_size=0.2)

dqn = DQN()
dqn.fit(train_episodes,
        eval_episodes=test_episodes,
        n_epochs=1,
        scorers={'td_error': td_error_scorer})
```

3.12.2 cross_validate

cross validation is also easily performed.

```
from d3rlpy.algos import DQN
from d3rlpy.datasets import get_cartpole
from d3rlpy.metrics import td_error_scorer
from sklearn.model_selection import cross_validate

dataset, env = get_cartpole()

dqn = DQN()
```

(continues on next page)

(continued from previous page)

```
scores = cross_validate(dqn,
                        dataset,
                        scoring={'td_error': td_error_scorer},
                        fit_params={'n_epochs': 1})
```

3.12.3 GridSearchCV

You can also perform grid search to find good hyperparameters.

```
from d3rlpy.algos import DQN
from d3rlpy.datasets import get_cartpole
from d3rlpy.metrics import td_error_scorer
from sklearn.model_selection import GridSearchCV

dataset, env = get_cartpole()

dqn = DQN()

gscv = GridSearchCV(estimator=dqn,
                    param_grid={'learning_rate': [1e-4, 3e-4, 1e-3]},
                    scoring={'td_error': td_error_scorer},
                    refit=False)

gscv.fit(dataset.episodes, n_epochs=1)
```

3.12.4 parallel execution with multiple GPUs

Some scikit-learn utilities provide *n_jobs* option, which enable fitting process to run in parallel to boost productivity. Ideally, if you have multiple GPUs, the multiple processes use different GPUs for computational efficiency.

d3rlpy provides special device assignment mechanism to realize this.

```
from d3rlpy.algos import DQN
from d3rlpy.datasets import get_cartpole
from d3rlpy.metrics import td_error_scorer
from d3rlpy.context import parallel
from sklearn.model_selection import cross_validate

dataset, env = get_cartpole()

# enable GPU
dqn = DQN(use_gpu=True)

# automatically assign different GPUs for the 4 processes.
with parallel():
    scores = cross_validate(dqn,
                            dataset,
                            scoring={'td_error': td_error_scorer},
                            fit_params={'n_epochs': 1},
                            n_jobs=4)
```

If *use_gpu=True* is passed, d3rlpy internally manages GPU device id via `d3rlpy.gpu.Device` object. This object is designed for scikit-learn's multi-process implementation that makes deep copies of the estimator object before dispatching. The *Device* object will increment its device id when deeply copied under the parallel context.

```

import copy
from d3rlpy.context import parallel
from d3rlpy.gpu import Device

device = Device(0)
# device.get_id() == 0

new_device = copy.deepcopy(device)
# new_device.get_id() == 0

with parallel():
    new_device = copy.deepcopy(device)
    # new_device.get_id() == 1
    # device.get_id() == 1

    new_device = copy.deepcopy(device)
    # if you have only 2 GPUs, it goes back to 0.
    # new_device.get_id() == 0
    # device.get_id() == 0

from d3rlpy.algos import DQN

dqn = DQN(use_gpu=Device(0)) # assign id=0
dqn = DQN(use_gpu=Device(1)) # assign id=1

```

3.13 Online Training

3.13.1 Standard Training

d3rlpy provides not only offline training, but also online training utilities. Despite being designed for offline training algorithms, d3rlpy is flexible enough to be trained in an online manner with a few more utilities.

```

import gym

from d3rlpy.algos import DQN
from d3rlpy.online.buffers import ReplayBuffer
from d3rlpy.online.explorers import LinearDecayEpsilonGreedy

# setup environment
env = gym.make('CartPole-v0')
eval_env = gym.make('CartPole-v0')

# setup algorithm
dqn = DQN(batch_size=32,
          learning_rate=2.5e-4,
          target_update_interval=100,
          use_gpu=True)

# setup replay buffer
buffer = ReplayBuffer(maxlen=1000000, env=env)

# setup explorers
explorer = LinearDecayEpsilonGreedy(start_epsilon=1.0,
                                    end_epsilon=0.1,

```

(continues on next page)

(continued from previous page)

```

                                duration=10000)

# start training
dqn.fit_online(env,
               buffer,
               explorer=explorer, # you don't need this with probablistic policy_
↳ algorithms
               eval_env=eval_env,
               n_epochs=30,
               n_steps_per_epoch=1000,
               n_updates_per_epoch=100)

```

Replay Buffer

<code>d3rlpy.online.buffers.ReplayBuffer</code>	Standard Replay Buffer.
---	-------------------------

d3rlpy.online.buffers.ReplayBuffer

class d3rlpy.online.buffers.**ReplayBuffer**(*maxlen*, *env=None*, *episodes=None*, *create_mask=False*, *mask_size=1*)

Standard Replay Buffer.

Parameters

- **maxlen** (*int*) – the maximum number of data length.
- **env** (*gym.Env*) – gym-like environment to extract shape information.
- **episodes** (*list* (*d3rlpy.dataset.Episode*)) – list of episodes to initialize buffer.
- **create_mask** (*bool*) – flag to create bootstrapping mask.
- **mask_size** (*int*) – ensemble size for binary mask.

Methods

`__len__()`

Return type *int*

append (*observation*, *action*, *reward*, *terminal*, *clip_episode=None*)

Append observation, action, reward and terminal flag to buffer.

If the terminal flag is True, Monte-Carlo returns will be computed with an entire episode and the whole transitions will be appended.

Parameters

- **observation** (*numpy.ndarray*) – observation.
- **action** (*numpy.ndarray*) – action.
- **reward** (*float*) – reward.
- **terminal** (*float*) – terminal flag.

- **clip_episode** (*Optional[bool]*) – flag to clip the current episode. If None, the episode is clipped based on terminal.

Return type `None`

append_episode (*episode*)

Append Episode object to buffer.

Parameters **episode** (`d3rlpy.dataset.Episode`) – episode.

Return type `None`

sample (*batch_size, n_frames=1, n_steps=1, gamma=0.99*)

Returns sampled mini-batch of transitions.

If observation is image, you can stack arbitrary frames via `n_frames`.

```
buffer.observation_shape == (3, 84, 84)

# stack 4 frames
batch = buffer.sample(batch_size=32, n_frames=4)

batch.observations.shape == (32, 12, 84, 84)
```

Parameters

- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – the number of steps before the next observation.
- **gamma** (*float*) – discount factor used in N-step return calculation.

Returns mini-batch.

Return type `d3rlpy.dataset.TransitionMiniBatch`

size ()

Returns the number of appended elements in buffer.

Returns the number of elements in buffer.

Return type `int`

to_mdp_dataset ()

Convert replay data into static dataset.

The length of the dataset can be longer than the length of the replay buffer because this conversion is done by tracing `Transition` objects.

Returns `MDPDataSet` object.

Return type `d3rlpy.dataset.MDPDataSet`

Attributes

transitions

Returns a FIFO queue of transitions.

Returns FIFO queue of transitions.

Return type d3rlpy.online.buffers.FIFOQueue

Explorers

<code>d3rlpy.online.explorers.ConstantEpsilonGreedy</code>	ϵ -greedy explorer with constant ϵ .
<code>d3rlpy.online.explorers.LinearDecayEpsilonGreedy</code>	ϵ -greedy explorer with linear decay schedule.
<code>d3rlpy.online.explorers.NormalNoise</code>	Normal noise explorer.

d3rlpy.online.explorers.ConstantEpsilonGreedy

class d3rlpy.online.explorers.ConstantEpsilonGreedy(*epsilon*)

ϵ -greedy explorer with constant ϵ .

Parameters *epsilon* (*float*) – the constant ϵ .

Methods

sample (*algo*, *x*, *step*)

Parameters

- **algo** (d3rlpy.online.explorers._ActionProtocol) –
- **x** (numpy.ndarray) –
- **step** (int) –

Return type numpy.ndarray

d3rlpy.online.explorers.LinearDecayEpsilonGreedy

class d3rlpy.online.explorers.LinearDecayEpsilonGreedy(*start_epsilon=1.0*,
end_epsilon=0.1, *duration=1000000*
end_epsilon=0.1, *duration=1000000*)

ϵ -greedy explorer with linear decay schedule.

Parameters

- **start_epsilon** (*float*) – the beginning ϵ .
- **end_epsilon** (*float*) – the end ϵ .
- **duration** (*int*) – the scheduling duration.

Methods

compute_epsilon (*step*)

Returns decayed ϵ .

Returns ϵ .

Parameters **step** (*int*) –

Return type *float*

sample (*algo*, *x*, *step*)

Returns ϵ -greedy action.

Parameters

- **algo** (*d3rlpy.online.explorers._ActionProtocol*) – algorithm.
- **x** (*numpy.ndarray*) – observation.
- **step** (*int*) – current environment step.

Returns ϵ -greedy action.

Return type *numpy.ndarray*

d3rlpy.online.explorers.NormalNoise

class *d3rlpy.online.explorers.NormalNoise* (*mean=0.0*, *std=0.1*)

Normal noise explorer.

Parameters

- **mean** (*float*) – mean.
- **std** (*float*) – standard deviation.

Methods

sample (*algo*, *x*, *step*)

Returns action with noise injection.

Parameters

- **algo** (*d3rlpy.online.explorers._ActionProtocol*) – algorithm.
- **x** (*numpy.ndarray*) – observation.
- **step** (*int*) –

Returns action with noise injection.

Return type *numpy.ndarray*

3.13.2 Batch Concurrent Training

d3rlpy supports computationally efficient batch concurrent training.

```
import gym

from d3rlpy.algos import DQN
from d3rlpy.envs import AsyncBatchEnv
from d3rlpy.online.buffers import BatchReplayBuffer
from d3rlpy.online.explorers import LinearDecayEpsilonGreedy

# this condition is necessary due to spawning processes
if __name__ == '__main__':
    env = AsyncBatchEnv([lambda: gym.make('CartPole-v0') for _ in range(10)])

    eval_env = gym.make('CartPole-v0')

    # setup algorithm
    dqn = DQN(batch_size=32,
              learning_rate=2.5e-4,
              target_update_interval=100,
              use_gpu=True)

    # setup replay buffer
    buffer = BatchReplayBuffer(maxlen=1000000, env=env)

    # setup explorers
    explorer = LinearDecayEpsilonGreedy(start_epsilon=1.0,
                                       end_epsilon=0.1,
                                       duration=10000)

    # start training
    dqn.fit_batch_online(env,
                        buffer,
                        explorer=explorer, # you don't need this with probabilistic_
    ↪ policy algorithms
                        eval_env=eval_env,
                        n_epochs=30,
                        n_steps_per_epoch=1000,
                        n_updates_per_epoch=100)
```

For the environment wrapper, please see `d3rlpy.envs.AsyncBatchEnv` and `d3rlpy.envs.SyncBatchEnv`.

Replay Buffer

`d3rlpy.online.buffers.`
`BatchReplayBuffer`

Standard Replay Buffer for batch training.

d3rlpy.online.buffers.BatchReplayBuffer

class d3rlpy.online.buffers.**BatchReplayBuffer**(*maxlen*, *env*, *episodes=None*, *create_mask=False*, *mask_size=1*)

Standard Replay Buffer for batch training.

Parameters

- **maxlen** (*int*) – the maximum number of data length.
- **n_envs** (*int*) – the number of environments.
- **env** (*gym.Env*) – gym-like environment to extract shape information.
- **episodes** (*list* (*d3rlpy.dataset.Episode*)) – list of episodes to initialize buffer
- **create_mask** (*bool*) – flag to create bootstrapping mask.
- **mask_size** (*int*) – ensemble size for binary mask.

Methods

__len__ ()

Return type *int*

append (*observations*, *actions*, *rewards*, *terminals*, *clip_episodes=None*)

Append observation, action, reward and terminal flag to buffer.

If the terminal flag is True, Monte-Carlo returns will be computed with an entire episode and the whole transitions will be appended.

Parameters

- **observations** (*numpy.ndarray*) – observation.
- **actions** (*numpy.ndarray*) – action.
- **rewards** (*numpy.ndarray*) – reward.
- **terminals** (*numpy.ndarray*) – terminal flag.
- **clip_episodes** (*Optional* [*numpy.ndarray*]) – flag to clip the current episode.
If None, the episode is clipped based on `terminal`.

Return type *None*

append_episode (*episode*)

Append Episode object to buffer.

Parameters **episode** (*d3rlpy.dataset.Episode*) – episode.

Return type *None*

sample (*batch_size*, *n_frames=1*, *n_steps=1*, *gamma=0.99*)

Returns sampled mini-batch of transitions.

If observation is image, you can stack arbitrary frames via `n_frames`.

```
buffer.observation_shape == (3, 84, 84)

# stack 4 frames
batch = buffer.sample(batch_size=32, n_frames=4)
```

(continues on next page)

(continued from previous page)

```
batch.observations.shape == (32, 12, 84, 84)
```

Parameters

- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_steps** (*int*) – the number of steps before the next observation.
- **gamma** (*float*) – discount factor used in N-step return calculation.

Returns mini-batch.**Return type** *d3rlpy.dataset.TransitionMiniBatch***size()**

Returns the number of appended elements in buffer.

Returns the number of elements in buffer.**Return type** *int***to_mdp_dataset()**

Convert replay data into static dataset.

The length of the dataset can be longer than the length of the replay buffer because this conversion is done by tracing *Transition* objects.

Returns *MDPDataSet* object.**Return type** *d3rlpy.dataset.MDPDataSet***Attributes****transitions**

Returns a FIFO queue of transitions.

Returns FIFO queue of transitions.**Return type** *d3rlpy.online.buffers.FIFOQueue*

3.14 Model-based Algorithms

d3rlpy provides model-based reinforcement learning algorithms.

```
from d3rlpy.datasets import get_pendulum
from d3rlpy.dynamics import ProbabilisticEnsembleDynamics
from d3rlpy.metrics.scorer import dynamics_observation_prediction_error_scorer
from d3rlpy.metrics.scorer import dynamics_reward_prediction_error_scorer
from d3rlpy.metrics.scorer import dynamics_prediction_variance_scorer
from sklearn.model_selection import train_test_split

dataset, _ = get_pendulum()

train_episodes, test_episodes = train_test_split(dataset)
```

(continues on next page)

(continued from previous page)

```

dynamics = d3rlpy.dynamics.ProbabilisticEnsembleDynamics(learning_rate=1e-4, use_
    ↪ gpu=True)

# same as algorithms
dynamics.fit(train_episodes,
             eval_episodes=test_episodes,
             n_epochs=100,
             scorers={
                 'observation_error': dynamics_observation_prediction_error_scorer,
                 'reward_error': dynamics_reward_prediction_error_scorer,
                 'variance': dynamics_prediction_variance_scorer,
             })

```

Pick the best model and pass it to the model-based RL algorithm.

```

from d3rlpy.algos import MOPO

# load trained dynamics model
dynamics = ProbabilisticEnsembleDynamics.from_json('<path-to-params.json>/params.json
    ↪ ')
dynamics.load_model('<path-to-model>/model_xx.pt')

# give mopo as generator argument.
mopo = MOPO(dynamics=dynamics)

```

3.14.1 Dynamics Model

`d3rlpy.dynamics.`
`ProbabilisticEnsembleDynamics`

Probabilistic ensemble dynamics.

`d3rlpy.dynamics.ProbabilisticEnsembleDynamics`

```

class d3rlpy.dynamics.ProbabilisticEnsembleDynamics(*, learning_rate=0.001, op-
    tim_factory=d3rlpy.models.optimizers.AdamFactory(opti-
    betas=(0.9, 0.999), eps=1e-
    08, weight_decay=0.0001,
    amsgrad=False), en-
    coder_factory='default',
    batch_size=100, n_frames=1,
    n_ensembles=5, vari-
    ance_type='max', dis-
    crete_action=False,
    scaler=None, ac-
    tion_scaler=None,
    use_gpu=False, impl=None,
    **kwargs)

```

Probabilistic ensemble dynamics.

The ensemble dynamics model consists of N probabilistic models $\{T_{\theta_i}\}_{i=1}^N$. At each epoch, new transitions are

generated via randomly picked dynamics model T_θ .

$$s_{t+1}, r_{t+1} \sim T_\theta(s_t, a_t)$$

where $s_t \sim D$ for the first step, otherwise s_t is the previous generated observation, and $a_t \sim \pi(\cdot|s_t)$.

Note: Currently, `ProbabilisticEnsembleDynamics` only supports vector observations.

References

- Yu et al., [MOPO: Model-based Offline Policy Optimization](#).

Parameters

- **learning_rate** (*float*) – learning rate for dynamics model.
- **optim_factory** (`d3rlpy.models.optimizers.OptimizerFactory`) – optimizer factory.
- **encoder_factory** (`d3rlpy.models.encoders.EncoderFactory` or *str*) – encoder factory.
- **batch_size** (*int*) – mini-batch size.
- **n_frames** (*int*) – the number of frames to stack for image observation.
- **n_ensembles** (*int*) – the number of dynamics model for ensemble.
- **variance_type** (*str*) – variance calculation type. The available options are ['max', 'data'].
- **discrete_action** (*bool*) – flag to take discrete actions.
- **scaler** (`d3rlpy.preprocessing.scalers.Scaler` or *str*) – preprocessor. The available options are ['pixel', 'min_max', 'standard'].
- **action_scaler** (`d3rlpy.preprocessing.Actionscalers` or *str*) – action preprocessor. The available options are ['min_max'].
- **use_gpu** (*bool* or `d3rlpy.gpu.Device`) – flag to use GPU or device.
- **impl** (`d3rlpy.dynamics.torch.ProbabilisticEnsembleDynamicsImpl`) – dynamics implementation.

Methods

build_with_dataset (*dataset*)

Instantiate implementation object with `MDPDataSet` object.

Parameters **dataset** (`d3rlpy.dataset.MDPDataSet`) – dataset.

Return type `None`

build_with_env (*env*)

Instantiate implementation object with OpenAI Gym object.

Parameters **env** (`gym.core.Env`) – gym-like environment.

Return type `None`

create_impl (*observation_shape*, *action_size*)

Instantiate implementation objects with the dataset shapes.

This method will be used internally when *fit* method is called.

Parameters

- **observation_shape** (*Sequence[int]*) – observation shape.
- **action_size** (*int*) – dimension of action-space.

Return type `None`

fit (*dataset*, *n_epochs*=None, *n_steps*=None, *n_steps_per_epoch*=10000, *save_metrics*=True, *experiment_name*=None, *with_timestamp*=True, *logdir*='d3rlpy_logs', *verbose*=True, *show_progress*=True, *tensorboard_dir*=None, *eval_episodes*=None, *save_interval*=1, *scorers*=None, *shuffle*=True, *callback*=None)

Trains with the given dataset.

```
algo.fit(episodes, n_steps=1000000)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is None.
- **save_metrics** (*bool*) – flag to record metrics in files. If False, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if None, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (algo, epoch, total_step), which is called at the end of epochs.

Returns list of result tuples (epoch, metrics) per epoch.

Return type List[Tuple[int, Dict[str, float]]]

fitter (*dataset*, *n_epochs=None*, *n_steps=None*, *n_steps_per_epoch=10000*, *save_metrics=True*, *experiment_name=None*, *with_timestamp=True*, *logdir='d3rlpy_logs'*, *verbose=True*, *show_progress=True*, *tensorboard_dir=None*, *eval_episodes=None*, *save_interval=1*, *scorers=None*, *shuffle=True*, *callback=None*)

Iterate over epochs steps to train with the given dataset. At each iteration algo methods and properties can be changed or queried.

```
for epoch, metrics in algo.fitter(epochs):
    my_plot(metrics)
    algo.save_model(my_path)
```

Parameters

- **dataset** (*Union[List[d3rlpy.dataset.Episode], d3rlpy.dataset.MDPDataset]*) – list of episodes to train.
- **n_epochs** (*Optional[int]*) – the number of epochs to train.
- **n_steps** (*Optional[int]*) – the number of steps to train.
- **n_steps_per_epoch** (*int*) – the number of steps per epoch. This value will be ignored when *n_steps* is *None*.
- **save_metrics** (*bool*) – flag to record metrics in files. If *False*, the log directory is not created and the model parameters are not saved during training.
- **experiment_name** (*Optional[str]*) – experiment name for logging. If not passed, the directory name will be *{class name}_{timestamp}*.
- **with_timestamp** (*bool*) – flag to add timestamp string to the last of directory name.
- **logdir** (*str*) – root directory name to save logs.
- **verbose** (*bool*) – flag to show logged information on stdout.
- **show_progress** (*bool*) – flag to show progress bar for iterations.
- **tensorboard_dir** (*Optional[str]*) – directory to save logged information in tensorboard (additional to the csv data). if *None*, the directory will not be created.
- **eval_episodes** (*Optional[List[d3rlpy.dataset.Episode]]*) – list of episodes to test.
- **save_interval** (*int*) – interval to save parameters.
- **scorers** (*Optional[Dict[str, Callable[[Any, List[d3rlpy.dataset.Episode]], float]]]*) – list of scorer functions used with *eval_episodes*.
- **shuffle** (*bool*) – flag to shuffle transitions on each epoch.
- **callback** (*Optional[Callable[[d3rlpy.base.LearnableBase, int, int], None]]*) – callable function that takes (*algo*, *epoch*, *total_step*), which is called at the end of epochs.

Returns iterator yielding current epoch and metrics dict.

Return type Generator[Tuple[int, Dict[str, float]], None, None]

classmethod `from_json` (*fname*, *use_gpu=False*)

Returns algorithm configured with json file.

The Json file should be the one saved during fitting.

```
from d3rlpy.algos import Algo

# create algorithm with saved configuration
algo = Algo.from_json('d3rlpy_logs/<path-to-json>/params.json')

# ready to load
algo.load_model('d3rlpy_logs/<path-to-model>/model_100.pt')

# ready to predict
algo.predict(...)
```

Parameters

- **fname** (*str*) – file path to *params.json*.
- **use_gpu** (*Optional[Union[bool, int, d3rlpy.gpu.Device]]*) – flag to use GPU, device ID or device.

Returns algorithm.

Return type `d3rlpy.base.LearnableBase`

generate_new_data (*epoch*, *total_step*, *transitions*)

Returns generated transitions for data augmentation.

This method is for model-based RL algorithms.

Parameters

- **epoch** (*int*) – the current epoch.
- **total_step** (*int*) – the total update steps.
- **transitions** (*List[d3rlpy.dataset.Transition]*) – list of transitions.

Returns list of new transitions.

Return type `Optional[List[d3rlpy.dataset.Transition]]`

get_action_type ()

Returns action type (continuous or discrete).

Returns action type.

Return type `d3rlpy.constants.ActionSpace`

get_params (*deep=True*)

Returns the all attributes.

This method returns the all attributes including ones in subclasses. Some of scikit-learn utilities will use this method.

```
params = algo.get_params(deep=True)

# the returned values can be used to instantiate the new object.
algo2 = AlgoBase(**params)
```

Parameters **deep** (*bool*) – flag to deeply copy objects such as *impl*.

Returns attribute values in dictionary.

Return type Dict[str, Any]

load_model (*fname*)

Load neural network parameters.

```
algo.load_model('model.pt')
```

Parameters **fname** (*str*) – source file path.

Return type None

predict (*x*, *action*, *with_variance=False*, *indices=None*)

Returns predicted observation and reward.

Parameters

- **x** (*Union[numpy.ndarray, List[Any]]*) – observation
- **action** (*Union[numpy.ndarray, List[Any]]*) – action
- **with_variance** (*bool*) – flag to return prediction variance.
- **indices** (*Optional[numpy.ndarray]*) – index of ensemble model to return.

Returns tuple of predicted observation and reward. If *with_variance* is True, the prediction variance will be added as the 3rd element.

Return type Union[Tuple[numpy.ndarray, numpy.ndarray], Tuple[numpy.ndarray, numpy.ndarray, numpy.ndarray]]

save_model (*fname*)

Saves neural network parameters.

```
algo.save_model('model.pt')
```

Parameters **fname** (*str*) – destination file path.

Return type None

save_params (*logger*)

Saves configurations as params.json.

Parameters **logger** (*d3rlpy.logger.D3RLPyLogger*) – logger object.

Return type None

set_params (***params*)

Sets the given arguments to the attributes if they exist.

This method sets the given values to the attributes including ones in subclasses. If the values that don't exist as attributes are passed, they are ignored. Some of scikit-learn utilities will use this method.

```
algo.set_params(batch_size=100)
```

Parameters **params** (*Any*) – arbitrary inputs to set as attributes.

Returns itself.

Return type d3rlpy.base.LearnableBase

update (*epoch*, *total_step*, *batch*)

Update parameters with mini-batch of data.

Parameters

- **epoch** (*int*) – the current number of epochs.
- **total_step** (*int*) – the current number of total iterations.
- **batch** (`d3rlpy.dataset.TransitionMiniBatch`) – mini-batch data.

Returns dictionary of metrics.

Return type Dict[str, float]

Attributes

action_scaler

Preprocessing action scaler.

Returns preprocessing action scaler.

Return type Optional[ActionScaler]

action_size

Action size.

Returns action size.

Return type Optional[int]

batch_size

Batch size to train.

Returns batch size.

Return type int

gamma

Discount factor.

Returns discount factor.

Return type float

impl

Implementation object.

Returns implementation object.

Return type Optional[ImplBase]

n_frames

Number of frames to stack.

This is only for image observation.

Returns number of frames to stack.

Return type int

n_steps

N-step TD backup.

Returns N-step TD backup.

Return type int

observation_shape

Observation shape.

Returns observation shape.

Return type Optional[Sequence[int]]

scaler

Preprocessing scaler.

Returns preprocessing scaler.

Return type Optional[Scaler]

3.15 Stable-Baselines3 Wrapper

d3rlpy provides a minimal wrapper to use [Stable-Baselines3 \(SB3\)](#) features, like utility helpers or SB3 algorithms to create datasets.

Note: This wrapper is far from complete, and only provide a minimal integration with SB3.

3.15.1 Convert SB3 replay buffer to d3rlpy dataset

A replay buffer from Stable-Baselines3 can be easily converted to a `d3rlpy.dataset.MDPDataset` using `to_mdp_dataset()` utility function.

```
import stable_baselines3 as sb3

from d3rlpy.algos import AWR
from d3rlpy.wrappers.sb3 import to_mdp_dataset

# Train an off-policy agent with SB3
model = sb3.SAC("MlpPolicy", "Pendulum-v0", learning_rate=1e-3, verbose=1)
model.learn(6000)

# Convert to d3rlpy MDPDataset
dataset = to_mdp_dataset(model.replay_buffer)
# The dataset can then be used to train a d3rlpy model
offline_model = AWR()
offline_model.fit(dataset.episodes, n_epochs=100)
```

3.15.2 Convert d3rlpy to use SB3 helpers

An agent from d3rlpy can be converted to use the SB3 interface (notably follow the interface of SB3 `predict()`). This allows to use SB3 helpers like `evaluate_policy`.

```
import gym
from stable_baselines3.common.evaluation import evaluate_policy

from d3rlpy.algos import AWAC
from d3rlpy.wrappers.sb3 import SB3Wrapper
```

(continues on next page)

(continued from previous page)

```
env = gym.make("Pendulum-v0")

# Define an offline RL model
offline_model = AWAC()
# Train it using for instance a dataset created by a SB3 agent (see above)
offline_model.fit(dataset.episodes, n_epochs=10)

# Use SB3 wrapper (convert `predict()` method to follow SB3 API)
# to have access to SB3 helpers
# d3rlpy model is accessible via `wrapped_model.algo`
wrapped_model = SB3Wrapper(offline_model)

observation = env.reset()

# We can now use SB3's predict style
# it returns the action and the hidden states (for RNN policies)
action, _ = wrapped_model.predict([observation], deterministic=True)
# The following is equivalent to offline_model.sample_action(obs)
action, _ = wrapped_model.predict([observation], deterministic=False)

# Evaluate the trained model using SB3 helper
mean_reward, std_reward = evaluate_policy(wrapped_model, env)

print(f"mean_reward={mean_reward} +/- {std_reward}")

# Call methods from the wrapped d3rlpy model
wrapped_model.sample_action([observation])
wrapped_model.fit(dataset.episodes, n_epochs=10)

# Set attributes
wrapped_model.n_epochs = 2
# wrapped_model.n_epochs points to d3rlpy wrapped_model.algo.n_epochs
assert wrapped_model.algo.n_epochs == 2
```


COMMAND LINE INTERFACE

d3rlpy provides the convenient CLI tool.

4.1 plot

Plot the saved metrics by specifying paths:

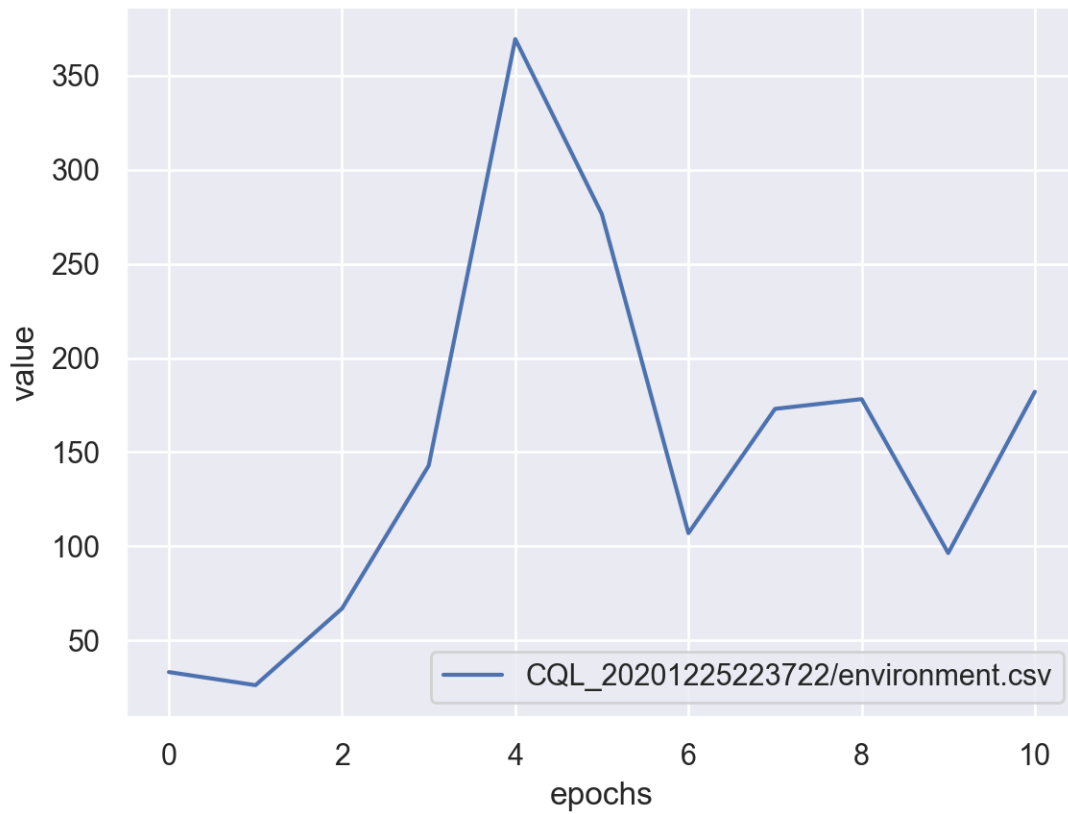
```
$ d3rlpy plot <path> [<path>...]
```

Table 1: options

option	description
--window	moving average window.
--show-steps	use iterations on x-axis.
--show-max	show maximum value.
--label	label in legend.
--xlim	limit on x-axis (tuple).
--ylim	limit on y-axis (tuple).
--title	title of the plot.

example:

```
$ d3rlpy plot d3rlpy_logs/CQL_20201224224314/environment.csv
```



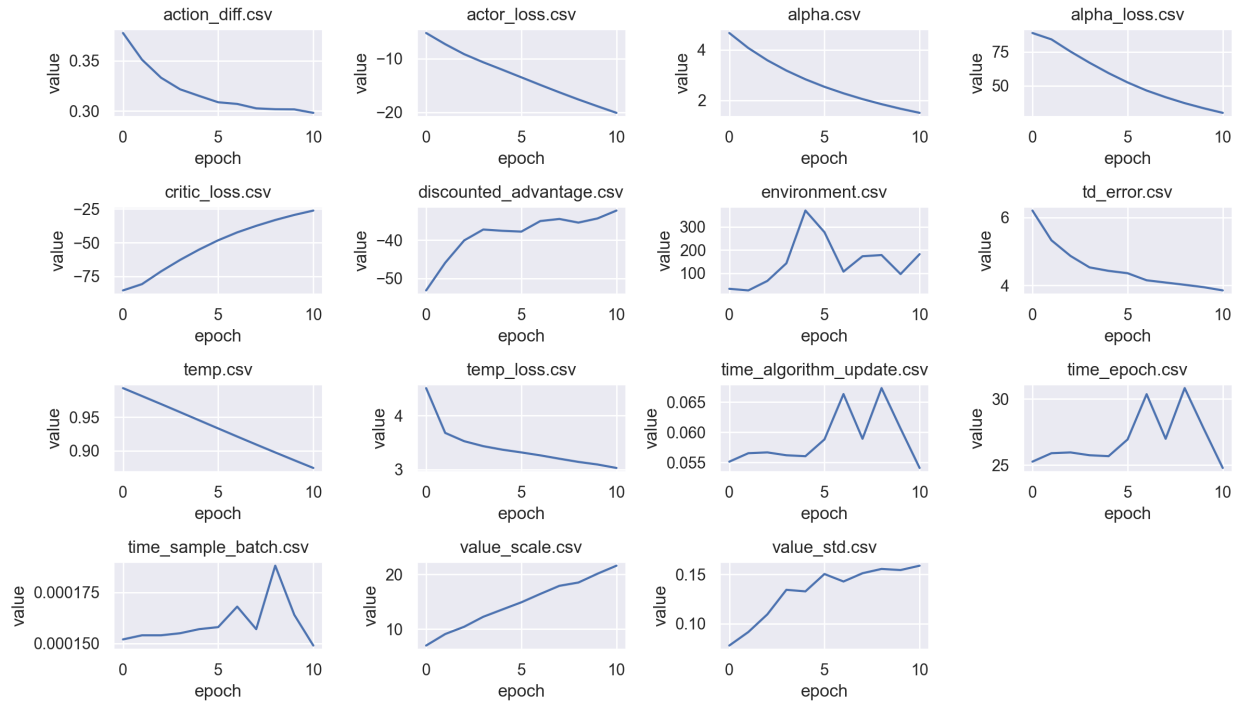
4.2 plot-all

Plot the all metrics saved in the directory:

```
$ d3rlpy plot-all <path>
```

example:

```
$ d3rlpy plot-all d3rlpy_logs/CQL_20201224224314
```



4.3 export

Export the saved model to the inference format, `onnx` and `torchscript`:

```
$ d3rlpy export <path>
```

Table 2: options

option	description
<code>--format</code>	model format (<code>torchscript</code> , <code>onnx</code>).
<code>--params-json</code>	explicitly specify <code>params.json</code> .
<code>--out</code>	output path.

example:

```
$ d3rlpy export d3rlpy_logs/CQL_20201224224314/model_100.pt
```

4.4 record

Record evaluation episodes as videos with the saved model:

```
$ d3rlpy record <path> --env-id <environment id>
```

Table 3: options

option	description
--env-id	Gym environment id.
--env-header	arbitrary Python code to define environment to evaluate.
--out	output directory.
--params-json	explicitly specify params.json
--n-episodes	the number of episodes to record.
--frame-rate	video frame rate.
--record-rate	images are recored every record-rate frames.
--epsilon	ϵ -greedy evaluation.

example:

```
# record simple environment
$ d3rlpy record d3rlpy_logs/CQL_20201224224314/model_100.pt --env-id HopperBulletEnv-
↪v0

# record wrapped environment
$ d3rlpy record d3rlpy_logs/Discrete_CQL_20201224224314/model_100.pt \
  --env-header 'import gym; from d3rlpy.envs import Atari; env = Atari(gym.make(
↪"BreakoutNoFrameskip-v4"), is_eval=True)'
```

4.5 play

Run evaluation episodes with rendering:

```
$ d3rlpy play <path> --env-id <environment id>
```

Table 4: options

option	description
--env-id	Gym environment id.
--env-header	arbitrary Python code to define environment to evaluate.
--params-json	explicitly specify params.json
--n-episodes	the number of episodes to run.

example:

```
# record simple environment
$ d3rlpy play d3rlpy_logs/CQL_20201224224314/model_100.pt --env-id HopperBulletEnv-v0

# record wrapped environment
$ d3rlpy play d3rlpy_logs/Discrete_CQL_20201224224314/model_100.pt \
  --env-header 'import gym; from d3rlpy.envs import Atari; env = Atari(gym.make(
↪"BreakoutNoFrameskip-v4"), is_eval=True)'
```

INSTALLATION

5.1 Recommended Platforms

d3rlpy supports Linux, macOS and also Windows.

5.2 Install d3rlpy

5.2.1 Install via PyPI

pip is a recommended way to install d3rlpy:

```
$ pip install d3rlpy
```

5.2.2 Install via Anaconda

d3rlpy is also available on *conda-forge*:

```
$ conda install -c conda-forge d3rlpy
```

5.2.3 Install via Docker

d3rlpy is also available on Docker Hub:

```
$ docker run -it --gpus all --name d3rlpy takuseno/d3rlpy:latest bash
```

5.2.4 Install from source

You can also install via GitHub repository:

```
$ git clone https://github.com/takuseno/d3rlpy
$ cd d3rlpy
$ pip install Cython numpy # if you have not installed them.
$ pip install -e .
```


6.1 Reproducibility

Reproducibility is one of the most important things when doing research activity. Here is a simple example in d3rlpy.

```
import d3rlpy
import gym

# set random seeds in random module, numpy module and PyTorch module.
d3rlpy.seed(313)

# set environment seed
env = gym.make('Hopper-v2')
env.seed(313)
```

6.2 Learning from image observation

d3rlpy supports both vector observations and image observations. There are several things you need to care about if you want to train RL agents from image observations.

```
from d3rlpy.dataset import MDPDataset

# observation MUST be uint8 array, and the channel-first images
observations = np.random.randint(256, size=(100000, 1, 84, 84), dtype=np.uint8)
actions = np.random.randint(4, size=100000)
rewards = np.random.random(100000)
terminals = np.random.randint(2, size=100000)

dataset = MDPDataset(observations, actions, rewards, terminals)

from d3rlpy.algos import DQN

dqn = DQN(scaler='pixel', # you MUST set pixel scaler
          n_frames=4) # you CAN set the number of frames to stack
```

6.3 Improve performance beyond the original paper

d3rlpy provides many options that you can use to improve performance potentially beyond the original paper. All the options are powerful, but the best combinations and hyperparameters are always dependent on the tasks.

```
from d3rlpy.models.encoders import DefaultEncoderFactory
from d3rlpy.models.q_functions import QRQFunctionFactory
from d3rlpy.algos import DQN, SAC

# use batch normalization
# this seems to improve performance with discrete action-spaces
encoder = DefaultEncoderFactory(use_batch_norm=True)

dqn = DQN(encoder_factory=encoder,
          n_critics=5, # Q function ensemble size
          n_steps=5, # N-step TD backup
          q_func_factory='qr') # use distributional Q function

# use dropout
# this will dramatically improve performance
encoder = DefaultEncoderFactory(dropout_rate=0.2)

sac = SAC(actor_encoder_factory=encoder)
```

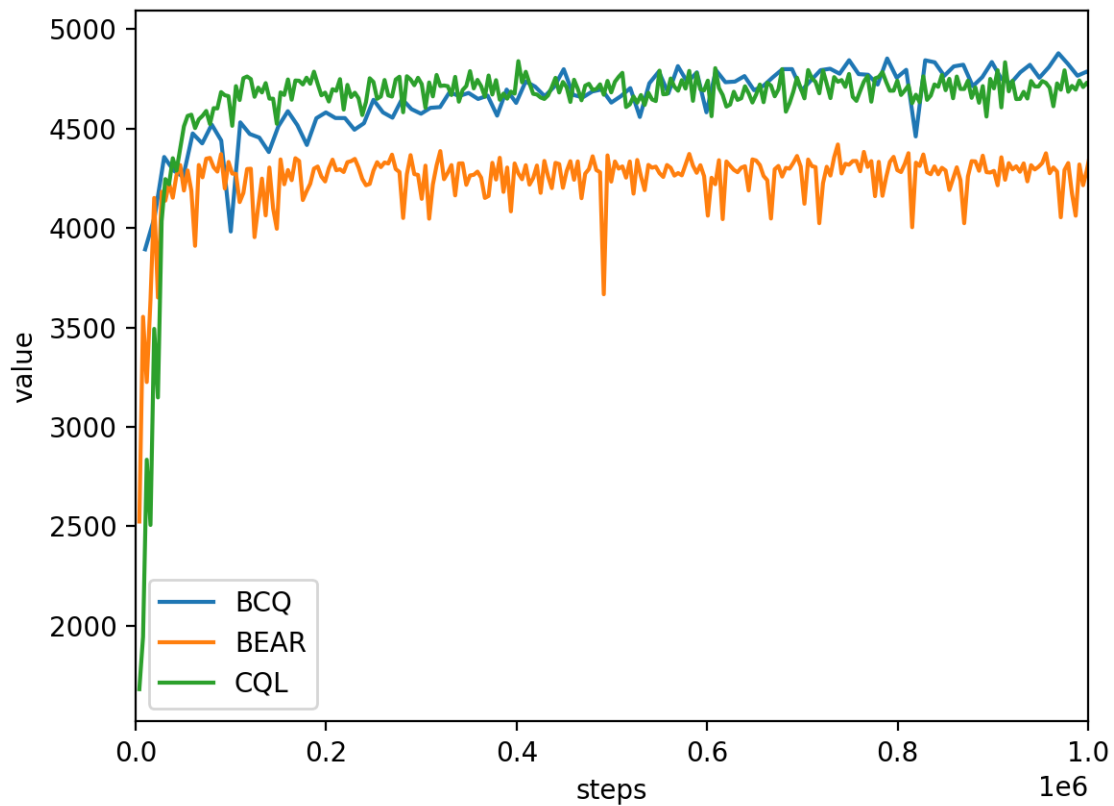

PAPER REPRODUCTIONS

For the experiment code, please take a look at [reproductions](#) directory.

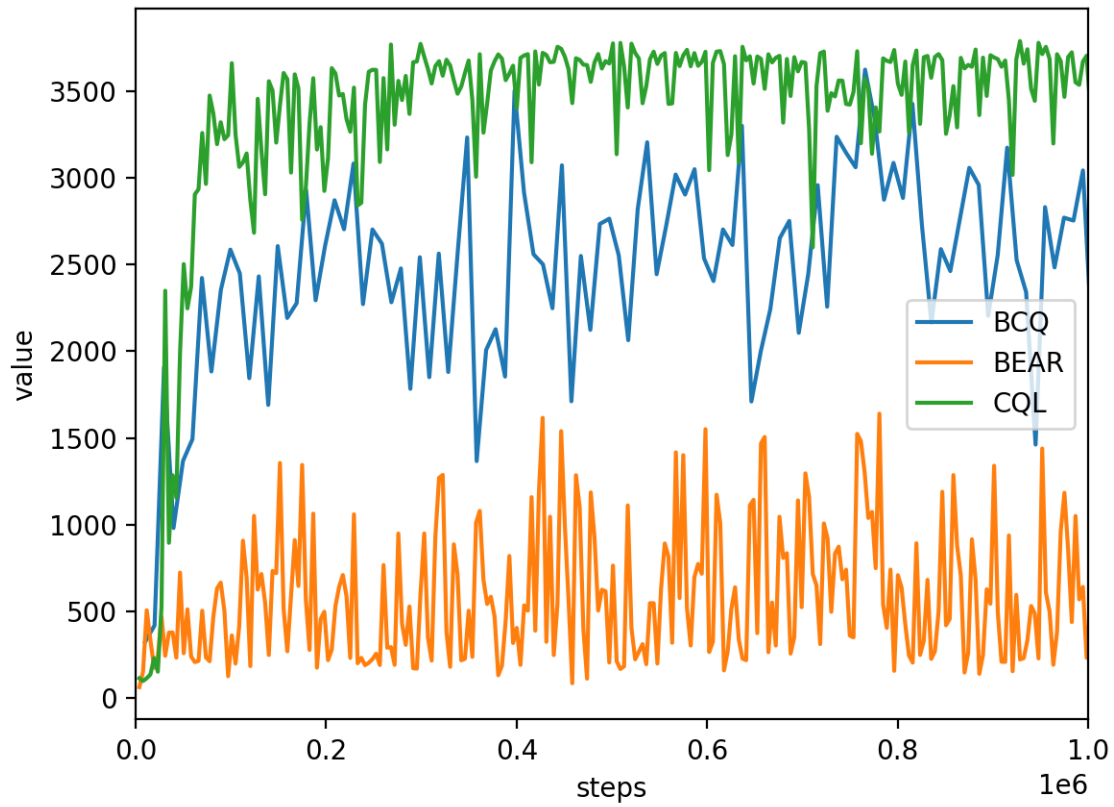
7.1 Offline

Fu et al., D4RL: Datasets for Deep Data-Driven Reinforcement Learning.

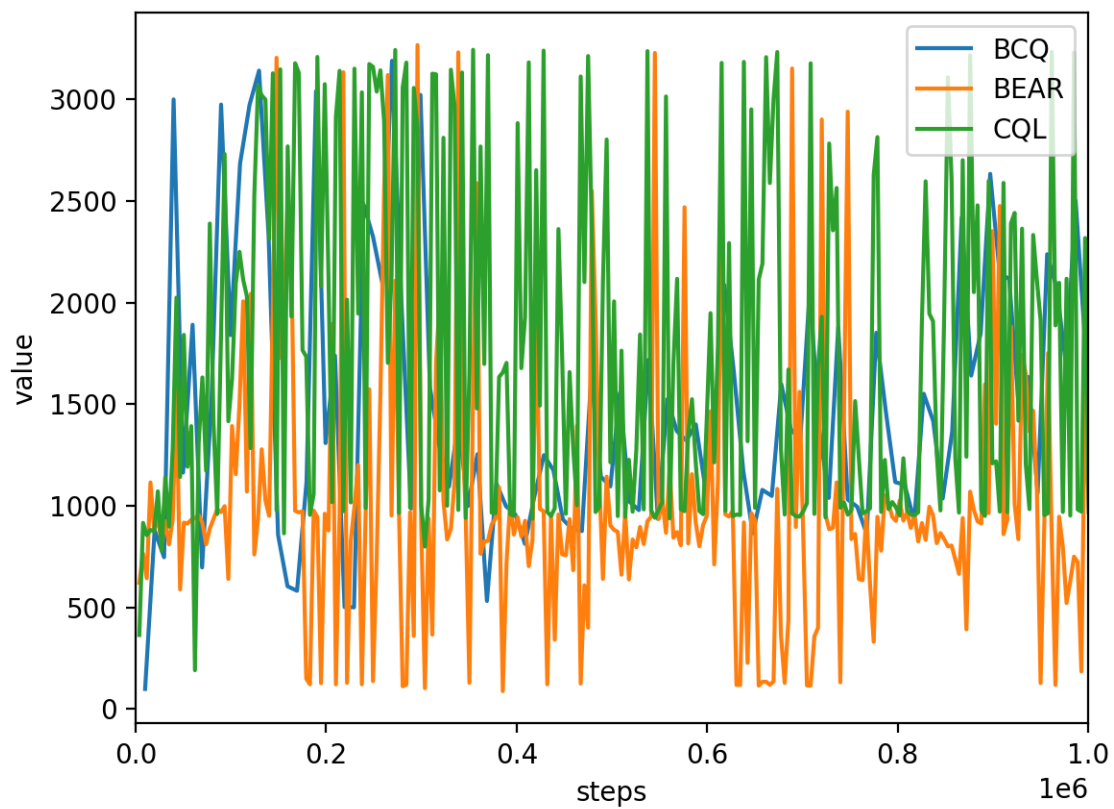
7.1.1 halfcheetah-medium-v0



7.1.2 walker2d-medium-v0



7.1.3 hopper-medium-v0



7.2 Online

TBD.

LICENSE

MIT License

Copyright (c) 2021 Takuma Seno

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

INDICES AND TABLES

- `genindex`
- `modindex`
- `search`

PYTHON MODULE INDEX

d

- `d3rlpy`, [9](#)
- `d3rlpy.algos`, [9](#)
- `d3rlpy.dataset`, [251](#)
- `d3rlpy.datasets`, [262](#)
- `d3rlpy.dynamics`, [323](#)
- `d3rlpy.metrics`, [283](#)
- `d3rlpy.models.encoders`, [276](#)
- `d3rlpy.models.optimizers`, [273](#)
- `d3rlpy.models.q_functions`, [245](#)
- `d3rlpy.online`, [316](#)
- `d3rlpy.ope`, [291](#)
- `d3rlpy.preprocessing`, [265](#)

Symbols

[__getitem__\(\)](#) (*d3rlpy.dataset.Episode method*), 256
[__getitem__\(\)](#) (*d3rlpy.dataset.MDPDataset method*), 253
[__getitem__\(\)](#) (*d3rlpy.dataset.TransitionMiniBatch method*), 260
[__iter__\(\)](#) (*d3rlpy.dataset.Episode method*), 256
[__iter__\(\)](#) (*d3rlpy.dataset.MDPDataset method*), 253
[__iter__\(\)](#) (*d3rlpy.dataset.TransitionMiniBatch method*), 260
[__len__\(\)](#) (*d3rlpy.dataset.Episode method*), 256
[__len__\(\)](#) (*d3rlpy.dataset.MDPDataset method*), 253
[__len__\(\)](#) (*d3rlpy.dataset.TransitionMiniBatch method*), 260
[__len__\(\)](#) (*d3rlpy.online.buffers.BatchReplayBuffer method*), 322
[__len__\(\)](#) (*d3rlpy.online.buffers.ReplayBuffer method*), 317

A

[action](#) (*d3rlpy.dataset.Transition attribute*), 259
[action_scaler](#) (*d3rlpy.algos.AWAC attribute*), 113
[action_scaler](#) (*d3rlpy.algos.AWR attribute*), 103
[action_scaler](#) (*d3rlpy.algos.BC attribute*), 18
[action_scaler](#) (*d3rlpy.algos.BCQ attribute*), 60
[action_scaler](#) (*d3rlpy.algos.BEAR attribute*), 71
[action_scaler](#) (*d3rlpy.algos.COMBO attribute*), 156
[action_scaler](#) (*d3rlpy.algos.CQL attribute*), 93
[action_scaler](#) (*d3rlpy.algos.CRR attribute*), 82
[action_scaler](#) (*d3rlpy.algos.DDPG attribute*), 28
[action_scaler](#) (*d3rlpy.algos.DiscreteAWR attribute*), 235
[action_scaler](#) (*d3rlpy.algos.DiscreteBC attribute*), 175
[action_scaler](#) (*d3rlpy.algos.DiscreteBCQ attribute*), 215
[action_scaler](#) (*d3rlpy.algos.DiscreteCQL attribute*), 225
[action_scaler](#) (*d3rlpy.algos.DiscreteRandomPolicy attribute*), 244

[action_scaler](#) (*d3rlpy.algos.DiscreteSAC attribute*), 205
[action_scaler](#) (*d3rlpy.algos.DoubleDQN attribute*), 195
[action_scaler](#) (*d3rlpy.algos.DQN attribute*), 185
[action_scaler](#) (*d3rlpy.algos.MOPO attribute*), 145
[action_scaler](#) (*d3rlpy.algos.PLAS attribute*), 123
[action_scaler](#) (*d3rlpy.algos.PLASWithPerturbation attribute*), 134
[action_scaler](#) (*d3rlpy.algos.RandomPolicy attribute*), 165
[action_scaler](#) (*d3rlpy.algos.SAC attribute*), 49
[action_scaler](#) (*d3rlpy.algos.TD3 attribute*), 38
[action_scaler](#) (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics attribute*), 330
[action_scaler](#) (*d3rlpy.ope.DiscreteFQE attribute*), 310
[action_scaler](#) (*d3rlpy.ope.FQE attribute*), 300
[action_size](#) (*d3rlpy.algos.AWAC attribute*), 113
[action_size](#) (*d3rlpy.algos.AWR attribute*), 103
[action_size](#) (*d3rlpy.algos.BC attribute*), 18
[action_size](#) (*d3rlpy.algos.BCQ attribute*), 60
[action_size](#) (*d3rlpy.algos.BEAR attribute*), 71
[action_size](#) (*d3rlpy.algos.COMBO attribute*), 156
[action_size](#) (*d3rlpy.algos.CQL attribute*), 93
[action_size](#) (*d3rlpy.algos.CRR attribute*), 82
[action_size](#) (*d3rlpy.algos.DDPG attribute*), 28
[action_size](#) (*d3rlpy.algos.DiscreteAWR attribute*), 235
[action_size](#) (*d3rlpy.algos.DiscreteBC attribute*), 175
[action_size](#) (*d3rlpy.algos.DiscreteBCQ attribute*), 215
[action_size](#) (*d3rlpy.algos.DiscreteCQL attribute*), 225
[action_size](#) (*d3rlpy.algos.DiscreteRandomPolicy attribute*), 244
[action_size](#) (*d3rlpy.algos.DiscreteSAC attribute*), 205
[action_size](#) (*d3rlpy.algos.DoubleDQN attribute*), 195
[action_size](#) (*d3rlpy.algos.DQN attribute*), 185
[action_size](#) (*d3rlpy.algos.MOPO attribute*), 145

- action_size (*d3rlpy.algos.PLAS* attribute), 123
 action_size (*d3rlpy.algos.PLASWithPerturbation* attribute), 134
 action_size (*d3rlpy.algos.RandomPolicy* attribute), 165
 action_size (*d3rlpy.algos.SAC* attribute), 49
 action_size (*d3rlpy.algos.TD3* attribute), 38
 action_size (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics* attribute), 330
 action_size (*d3rlpy.ope.DiscreteFQE* attribute), 310
 action_size (*d3rlpy.ope.FQE* attribute), 300
 actions (*d3rlpy.dataset.Episode* attribute), 257
 actions (*d3rlpy.dataset.MDPDataset* attribute), 255
 actions (*d3rlpy.dataset.TransitionMiniBatch* attribute), 261
 AdamFactory (class in *d3rlpy.models.optimizers*), 275
 add_additional_data() (*d3rlpy.dataset.TransitionMiniBatch* method), 260
 append() (*d3rlpy.dataset.MDPDataset* method), 253
 append() (*d3rlpy.online.buffer.BatchReplayBuffer* method), 322
 append() (*d3rlpy.online.buffer.ReplayBuffer* method), 317
 append_episode() (*d3rlpy.online.buffer.BatchReplayBuffer* method), 322
 append_episode() (*d3rlpy.online.buffer.ReplayBuffer* method), 318
 average_value_estimation_scorer() (in module *d3rlpy.metrics.scorer*), 285
 AWAC (class in *d3rlpy.algos*), 104
 AWR (class in *d3rlpy.algos*), 94
- ## B
- batch_size (*d3rlpy.algos.AWAC* attribute), 113
 batch_size (*d3rlpy.algos.AWR* attribute), 103
 batch_size (*d3rlpy.algos.BC* attribute), 18
 batch_size (*d3rlpy.algos.BCQ* attribute), 60
 batch_size (*d3rlpy.algos.BEAR* attribute), 71
 batch_size (*d3rlpy.algos.COMBO* attribute), 156
 batch_size (*d3rlpy.algos.CQL* attribute), 93
 batch_size (*d3rlpy.algos.CRR* attribute), 82
 batch_size (*d3rlpy.algos.DDPG* attribute), 28
 batch_size (*d3rlpy.algos.DiscreteAWR* attribute), 235
 batch_size (*d3rlpy.algos.DiscreteBC* attribute), 175
 batch_size (*d3rlpy.algos.DiscreteBCQ* attribute), 215
 batch_size (*d3rlpy.algos.DiscreteCQL* attribute), 225
 batch_size (*d3rlpy.algos.DiscreteRandomPolicy* attribute), 244
 batch_size (*d3rlpy.algos.DiscreteSAC* attribute), 205
 batch_size (*d3rlpy.algos.DoubleDQN* attribute), 195
 batch_size (*d3rlpy.algos.DQN* attribute), 185
 batch_size (*d3rlpy.algos.MOPO* attribute), 145
 batch_size (*d3rlpy.algos.PLAS* attribute), 123
 batch_size (*d3rlpy.algos.PLASWithPerturbation* attribute), 134
 batch_size (*d3rlpy.algos.RandomPolicy* attribute), 165
 batch_size (*d3rlpy.algos.SAC* attribute), 49
 batch_size (*d3rlpy.algos.TD3* attribute), 38
 batch_size (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics* attribute), 330
 batch_size (*d3rlpy.ope.DiscreteFQE* attribute), 310
 batch_size (*d3rlpy.ope.FQE* attribute), 300
 BatchReplayBuffer (class in *d3rlpy.online.buffer*), 322
 BC (class in *d3rlpy.algos*), 9
 BCQ (class in *d3rlpy.algos*), 50
 BEAR (class in *d3rlpy.algos*), 61
 bootstrap (*d3rlpy.models.q_functions.FQFQFunctionFactory* attribute), 250
 bootstrap (*d3rlpy.models.q_functions.IQNQFunctionFactory* attribute), 249
 bootstrap (*d3rlpy.models.q_functions.MeanQFunctionFactory* attribute), 247
 bootstrap (*d3rlpy.models.q_functions.QRQFunctionFactory* attribute), 248
 build_episodes() (*d3rlpy.dataset.MDPDataset* method), 253
 build_transitions() (*d3rlpy.dataset.Episode* method), 256
 build_with_dataset() (in module *d3rlpy.algos.AWAC* method), 105
 build_with_dataset() (*d3rlpy.algos.AWR* method), 95
 build_with_dataset() (*d3rlpy.algos.BC* method), 10
 build_with_dataset() (*d3rlpy.algos.BCQ* method), 52
 build_with_dataset() (*d3rlpy.algos.BEAR* method), 63
 build_with_dataset() (*d3rlpy.algos.COMBO* method), 148
 build_with_dataset() (*d3rlpy.algos.CQL* method), 85
 build_with_dataset() (*d3rlpy.algos.CRR* method), 74
 build_with_dataset() (*d3rlpy.algos.DDPG* method), 20
 build_with_dataset() (*d3rlpy.algos.DiscreteAWR* method), 227
 build_with_dataset() (*d3rlpy.algos.DiscreteBC* method), 167
 build_with_dataset() (*d3rlpy.algos.DiscreteBCQ* method), 207
 build_with_dataset() (*d3rlpy.algos.DiscreteCQL* method), 217

`build_with_dataset()`
 (*d3rlpy.algos.DiscreteRandomPolicy method*), 236
`build_with_dataset()` (*d3rlpy.algos.DiscreteSAC method*), 197
`build_with_dataset()` (*d3rlpy.algos.DoubleDQN method*), 187
`build_with_dataset()` (*d3rlpy.algos.DQN method*), 177
`build_with_dataset()` (*d3rlpy.algos.MOPO method*), 137
`build_with_dataset()` (*d3rlpy.algos.PLAS method*), 115
`build_with_dataset()`
 (*d3rlpy.algos.PLASWithPerturbation method*), 126
`build_with_dataset()`
 (*d3rlpy.algos.RandomPolicy method*), 157
`build_with_dataset()` (*d3rlpy.algos.SAC method*), 41
`build_with_dataset()` (*d3rlpy.algos.TD3 method*), 30
`build_with_dataset()`
 (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 325
`build_with_dataset()` (*d3rlpy.ope.DiscreteFQE method*), 302
`build_with_dataset()` (*d3rlpy.ope.FQE method*), 292
`build_with_env()` (*d3rlpy.algos.AWAC method*), 105
`build_with_env()` (*d3rlpy.algos.AWR method*), 95
`build_with_env()` (*d3rlpy.algos.BC method*), 10
`build_with_env()` (*d3rlpy.algos.BCQ method*), 52
`build_with_env()` (*d3rlpy.algos.BEAR method*), 63
`build_with_env()` (*d3rlpy.algos.COMBO method*), 148
`build_with_env()` (*d3rlpy.algos.CQL method*), 85
`build_with_env()` (*d3rlpy.algos.CRR method*), 74
`build_with_env()` (*d3rlpy.algos.DDPG method*), 20
`build_with_env()` (*d3rlpy.algos.DiscreteAWR method*), 227
`build_with_env()` (*d3rlpy.algos.DiscreteBC method*), 167
`build_with_env()` (*d3rlpy.algos.DiscreteBCQ method*), 207
`build_with_env()` (*d3rlpy.algos.DiscreteCQL method*), 217
`build_with_env()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 236
`build_with_env()` (*d3rlpy.algos.DiscreteSAC method*), 197
`build_with_env()` (*d3rlpy.algos.DoubleDQN method*), 187
`build_with_env()` (*d3rlpy.algos.DQN method*), 177
`build_with_env()` (*d3rlpy.algos.MOPO method*), 137
`build_with_env()` (*d3rlpy.algos.PLAS method*), 115
`build_with_env()` (*d3rlpy.algos.PLASWithPerturbation method*), 126
`build_with_env()` (*d3rlpy.algos.RandomPolicy method*), 157
`build_with_env()` (*d3rlpy.algos.SAC method*), 41
`build_with_env()` (*d3rlpy.algos.TD3 method*), 30
`build_with_env()` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 325
`build_with_env()` (*d3rlpy.ope.DiscreteFQE method*), 302
`build_with_env()` (*d3rlpy.ope.FQE method*), 292

C

`clear_links()` (*d3rlpy.dataset.Transition method*), 258
`clip_reward()` (*d3rlpy.dataset.MDPDataset method*), 253
`collect()` (*d3rlpy.algos.AWAC method*), 105
`collect()` (*d3rlpy.algos.AWR method*), 95
`collect()` (*d3rlpy.algos.BC method*), 10
`collect()` (*d3rlpy.algos.BCQ method*), 52
`collect()` (*d3rlpy.algos.BEAR method*), 63
`collect()` (*d3rlpy.algos.COMBO method*), 148
`collect()` (*d3rlpy.algos.CQL method*), 85
`collect()` (*d3rlpy.algos.CRR method*), 74
`collect()` (*d3rlpy.algos.DDPG method*), 20
`collect()` (*d3rlpy.algos.DiscreteAWR method*), 228
`collect()` (*d3rlpy.algos.DiscreteBC method*), 167
`collect()` (*d3rlpy.algos.DiscreteBCQ method*), 207
`collect()` (*d3rlpy.algos.DiscreteCQL method*), 217
`collect()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 236
`collect()` (*d3rlpy.algos.DiscreteSAC method*), 197
`collect()` (*d3rlpy.algos.DoubleDQN method*), 187
`collect()` (*d3rlpy.algos.DQN method*), 177
`collect()` (*d3rlpy.algos.MOPO method*), 137
`collect()` (*d3rlpy.algos.PLAS method*), 116
`collect()` (*d3rlpy.algos.PLASWithPerturbation method*), 126
`collect()` (*d3rlpy.algos.RandomPolicy method*), 157
`collect()` (*d3rlpy.algos.SAC method*), 41
`collect()` (*d3rlpy.algos.TD3 method*), 30
`collect()` (*d3rlpy.ope.DiscreteFQE method*), 302
`collect()` (*d3rlpy.ope.FQE method*), 292
COMBO (class in *d3rlpy.algos*), 146
`compare_continuous_action_diff()` (in module *d3rlpy.metrics.comparer*), 288

`compare_discrete_action_match()` (in module `d3rlpy.metrics.comparer`), 289
`compute_epsilon()` (`d3rlpy.online.explorers.LinearDecayEpsilonGreedy` method), 320
`compute_return()` (`d3rlpy.dataset.Episode` method), 256
`compute_stats()` (`d3rlpy.dataset.MDPDataset` method), 253
`ConstantEpsilonGreedy` (class in `d3rlpy.online.explorers`), 319
`continuous_action_diff_scorer()` (in module `d3rlpy.metrics.scorer`), 287
`CQL` (class in `d3rlpy.algos`), 83
`create()` (`d3rlpy.models.encoders.DefaultEncoderFactory` method), 279
`create()` (`d3rlpy.models.encoders.DenseEncoderFactory` method), 282
`create()` (`d3rlpy.models.encoders.PixelEncoderFactory` method), 280
`create()` (`d3rlpy.models.encoders.VectorEncoderFactory` method), 281
`create()` (`d3rlpy.models.optimizers.AdamFactory` method), 275
`create()` (`d3rlpy.models.optimizers.OptimizerFactory` method), 274
`create()` (`d3rlpy.models.optimizers.RMSpropFactory` method), 276
`create()` (`d3rlpy.models.optimizers.SGDFactory` method), 274
`create_continuous()` (`d3rlpy.models.q_functions.FQFQFunctionFactory` method), 250
`create_continuous()` (`d3rlpy.models.q_functions.IQNQFunctionFactory` method), 248
`create_continuous()` (`d3rlpy.models.q_functions.MeanQFunctionFactory` method), 246
`create_continuous()` (`d3rlpy.models.q_functions.QRQFunctionFactory` method), 247
`create_discrete()` (`d3rlpy.models.q_functions.FQFQFunctionFactory` method), 250
`create_discrete()` (`d3rlpy.models.q_functions.IQNQFunctionFactory` method), 248
`create_discrete()` (`d3rlpy.models.q_functions.MeanQFunctionFactory` method), 246
`create_discrete()` (`d3rlpy.models.q_functions.QRQFunctionFactory` method), 247
`create_impl()` (`d3rlpy.algos.AWAC` method), 106
`create_impl()` (`d3rlpy.algos.AWR` method), 95
`create_impl()` (`d3rlpy.algos.BC` method), 11
`create_impl()` (`d3rlpy.algos.BCQ` method), 52
`create_impl()` (`d3rlpy.algos.BEAR` method), 63
`create_impl()` (`d3rlpy.algos.COMBO` method), 148
`create_impl()` (`d3rlpy.algos.CQL` method), 85
`create_impl()` (`d3rlpy.algos.CRR` method), 74
`create_impl()` (`d3rlpy.algos.DDPG` method), 20
`create_impl()` (`d3rlpy.algos.DiscreteAWR` method), 228
`create_impl()` (`d3rlpy.algos.DiscreteBC` method), 167
`create_impl()` (`d3rlpy.algos.DiscreteBCQ` method), 208
`create_impl()` (`d3rlpy.algos.DiscreteCQL` method), 218
`create_impl()` (`d3rlpy.algos.DiscreteRandomPolicy` method), 236
`create_impl()` (`d3rlpy.algos.DiscreteSAC` method), 198
`create_impl()` (`d3rlpy.algos.DoubleDQN` method), 187
`create_impl()` (`d3rlpy.algos.DQN` method), 177
`create_impl()` (`d3rlpy.algos.MOPO` method), 137
`create_impl()` (`d3rlpy.algos.PLAS` method), 116
`create_impl()` (`d3rlpy.algos.PLASWithPerturbation` method), 127
`create_impl()` (`d3rlpy.algos.RandomPolicy` method), 157
`create_impl()` (`d3rlpy.algos.SAC` method), 41
`create_impl()` (`d3rlpy.algos.TD3` method), 31
`create_impl()` (`d3rlpy.dynamics.ProbabilisticEnsembleDynamics` method), 325
`create_impl()` (`d3rlpy.ope.DiscreteFQE` method), 303
`create_impl()` (`d3rlpy.ope.FQE` method), 293
`create_with_action()` (`d3rlpy.models.encoders.DefaultEncoderFactory` method), 279
`create_with_action()` (`d3rlpy.models.encoders.DenseEncoderFactory` method), 282
`create_with_action()` (`d3rlpy.models.encoders.PixelEncoderFactory` method), 280
`create_with_action()` (`d3rlpy.models.encoders.VectorEncoderFactory` method), 281
`CRR` (class in `d3rlpy.algos`), 72

D

`d3rlpy` module, 9

d3rlpy.algos
 module, 9
d3rlpy.dataset
 module, 251
d3rlpy.datasets
 module, 262
d3rlpy.dynamics
 module, 323
d3rlpy.metrics
 module, 283
d3rlpy.models.encoders
 module, 276
d3rlpy.models.optimizers
 module, 273
d3rlpy.models.q_functions
 module, 245
d3rlpy.online
 module, 316
d3rlpy.ope
 module, 291
d3rlpy.preprocessing
 module, 265
DDPG (class in d3rlpy.algos), 19
DefaultEncoderFactory (class in d3rlpy.models.encoders), 279
DenseEncoderFactory (class in d3rlpy.models.encoders), 282
discounted_sum_of_advantage_scorer() (in module d3rlpy.metrics.scorer), 285
discrete_action_match_scorer() (in module d3rlpy.metrics.scorer), 287
DiscreteAWR (class in d3rlpy.algos), 226
DiscreteBC (class in d3rlpy.algos), 166
DiscreteBCQ (class in d3rlpy.algos), 206
DiscreteCQL (class in d3rlpy.algos), 216
DiscreteFQE (class in d3rlpy.ope), 301
DiscreteRandomPolicy (class in d3rlpy.algos), 236
DiscreteSAC (class in d3rlpy.algos), 196
DoubleDQN (class in d3rlpy.algos), 186
DQN (class in d3rlpy.algos), 176
dump() (d3rlpy.dataset.MDPDataset method), 254
dynamics_observation_prediction_error_scorer() (in module d3rlpy.metrics.scorer), 290
dynamics_prediction_variance_scorer() (in module d3rlpy.metrics.scorer), 290
dynamics_reward_prediction_error_scorer() (in module d3rlpy.metrics.scorer), 290
entropy_coeff (d3rlpy.models.q_functions.FQFQFunctionFactory attribute), 250
Episode (class in d3rlpy.dataset), 256
episode_terminals (d3rlpy.dataset.MDPDataset attribute), 255
episodes (d3rlpy.dataset.MDPDataset attribute), 255
evaluate_on_environment() (in module d3rlpy.metrics.scorer), 288
extend() (d3rlpy.dataset.MDPDataset method), 254
F
fit() (d3rlpy.algos.AWAC method), 106
fit() (d3rlpy.algos.AWR method), 96
fit() (d3rlpy.algos.BC method), 11
fit() (d3rlpy.algos.BCQ method), 52
fit() (d3rlpy.algos.BEAR method), 64
fit() (d3rlpy.algos.COMBO method), 148
fit() (d3rlpy.algos.CQL method), 85
fit() (d3rlpy.algos.CRR method), 74
fit() (d3rlpy.algos.DDPG method), 21
fit() (d3rlpy.algos.DiscreteAWR method), 228
fit() (d3rlpy.algos.DiscreteBC method), 168
fit() (d3rlpy.algos.DiscreteBCQ method), 208
fit() (d3rlpy.algos.DiscreteCQL method), 218
fit() (d3rlpy.algos.DiscreteRandomPolicy method), 237
fit() (d3rlpy.algos.DiscreteSAC method), 198
fit() (d3rlpy.algos.DoubleDQN method), 187
fit() (d3rlpy.algos.DQN method), 177
fit() (d3rlpy.algos.MOPO method), 138
fit() (d3rlpy.algos.PLAS method), 116
fit() (d3rlpy.algos.PLASWithPerturbation method), 127
fit() (d3rlpy.algos.RandomPolicy method), 158
fit() (d3rlpy.algos.SAC method), 41
fit() (d3rlpy.algos.TD3 method), 31
fit() (d3rlpy.dynamics.ProbabilisticEnsembleDynamics method), 326
fit() (d3rlpy.ope.DiscreteFQE method), 303
fit() (d3rlpy.ope.FQE method), 293
fit() (d3rlpy.preprocessing.MinMaxActionScaler method), 272
fit() (d3rlpy.preprocessing.MinMaxScaler method), 268
fit() (d3rlpy.preprocessing.PixelScaler method), 266
fit() (d3rlpy.preprocessing.StandardScaler method), 269
fit_batch_online() (d3rlpy.algos.AWAC method), 107
fit_batch_online() (d3rlpy.algos.AWR method), 96
fit_batch_online() (d3rlpy.algos.BC method), 12
fit_batch_online() (d3rlpy.algos.BCQ method), 53
E
embed_size (d3rlpy.models.q_functions.FQFQFunctionFactory attribute), 250
embed_size (d3rlpy.models.q_functions.IQNQFunctionFactory attribute), 249

`fit_batch_online()` (*d3rlpy.algos.BEAR method*), 64
`fit_batch_online()` (*d3rlpy.algos.COMBO method*), 149
`fit_batch_online()` (*d3rlpy.algos.CQL method*), 86
`fit_batch_online()` (*d3rlpy.algos.CRR method*), 75
`fit_batch_online()` (*d3rlpy.algos.DDPG method*), 21
`fit_batch_online()` (*d3rlpy.algos.DiscreteAWR method*), 229
`fit_batch_online()` (*d3rlpy.algos.DiscreteBC method*), 168
`fit_batch_online()` (*d3rlpy.algos.DiscreteBCQ method*), 209
`fit_batch_online()` (*d3rlpy.algos.DiscreteCQL method*), 219
`fit_batch_online()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 238
`fit_batch_online()` (*d3rlpy.algos.DiscreteSAC method*), 199
`fit_batch_online()` (*d3rlpy.algos.DoubleDQN method*), 188
`fit_batch_online()` (*d3rlpy.algos.DQN method*), 178
`fit_batch_online()` (*d3rlpy.algos.MOPO method*), 138
`fit_batch_online()` (*d3rlpy.algos.PLAS method*), 117
`fit_batch_online()` (*d3rlpy.algos.PLASWithPerturbation method*), 128
`fit_batch_online()` (*d3rlpy.algos.RandomPolicy method*), 159
`fit_batch_online()` (*d3rlpy.algos.SAC method*), 42
`fit_batch_online()` (*d3rlpy.algos.TD3 method*), 32
`fit_batch_online()` (*d3rlpy.ope.DiscreteFQE method*), 304
`fit_batch_online()` (*d3rlpy.ope.FQE method*), 294
`fit_online()` (*d3rlpy.algos.AWAC method*), 108
`fit_online()` (*d3rlpy.algos.AWR method*), 97
`fit_online()` (*d3rlpy.algos.BC method*), 13
`fit_online()` (*d3rlpy.algos.BCQ method*), 54
`fit_online()` (*d3rlpy.algos.BEAR method*), 65
`fit_online()` (*d3rlpy.algos.COMBO method*), 150
`fit_online()` (*d3rlpy.algos.CQL method*), 87
`fit_online()` (*d3rlpy.algos.CRR method*), 76
`fit_online()` (*d3rlpy.algos.DDPG method*), 22
`fit_online()` (*d3rlpy.algos.DiscreteAWR method*), 230
`fit_online()` (*d3rlpy.algos.DiscreteBC method*), 169
`fit_online()` (*d3rlpy.algos.DiscreteBCQ method*), 210
`fit_online()` (*d3rlpy.algos.DiscreteCQL method*), 220
`fit_online()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 238
`fit_online()` (*d3rlpy.algos.DiscreteSAC method*), 200
`fit_online()` (*d3rlpy.algos.DoubleDQN method*), 189
`fit_online()` (*d3rlpy.algos.DQN method*), 179
`fit_online()` (*d3rlpy.algos.MOPO method*), 139
`fit_online()` (*d3rlpy.algos.PLAS method*), 118
`fit_online()` (*d3rlpy.algos.PLASWithPerturbation method*), 129
`fit_online()` (*d3rlpy.algos.RandomPolicy method*), 159
`fit_online()` (*d3rlpy.algos.SAC method*), 43
`fit_online()` (*d3rlpy.algos.TD3 method*), 33
`fit_online()` (*d3rlpy.ope.DiscreteFQE method*), 305
`fit_online()` (*d3rlpy.ope.FQE method*), 295
`fit_with_env()` (*d3rlpy.preprocessing.MinMaxActionScaler method*), 272
`fit_with_env()` (*d3rlpy.preprocessing.MinMaxScaler method*), 268
`fit_with_env()` (*d3rlpy.preprocessing.PixelScaler method*), 266
`fit_with_env()` (*d3rlpy.preprocessing.StandardScaler method*), 269
`fitter()` (*d3rlpy.algos.AWAC method*), 108
`fitter()` (*d3rlpy.algos.AWR method*), 98
`fitter()` (*d3rlpy.algos.BC method*), 14
`fitter()` (*d3rlpy.algos.BCQ method*), 55
`fitter()` (*d3rlpy.algos.BEAR method*), 66
`fitter()` (*d3rlpy.algos.COMBO method*), 151
`fitter()` (*d3rlpy.algos.CQL method*), 88
`fitter()` (*d3rlpy.algos.CRR method*), 77
`fitter()` (*d3rlpy.algos.DDPG method*), 23
`fitter()` (*d3rlpy.algos.DiscreteAWR method*), 231
`fitter()` (*d3rlpy.algos.DiscreteBC method*), 170
`fitter()` (*d3rlpy.algos.DiscreteBCQ method*), 211
`fitter()` (*d3rlpy.algos.DiscreteCQL method*), 221
`fitter()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 239
`fitter()` (*d3rlpy.algos.DiscreteSAC method*), 200
`fitter()` (*d3rlpy.algos.DoubleDQN method*), 190
`fitter()` (*d3rlpy.algos.DQN method*), 180
`fitter()` (*d3rlpy.algos.MOPO method*), 140
`fitter()` (*d3rlpy.algos.PLAS method*), 119
`fitter()` (*d3rlpy.algos.PLASWithPerturbation method*), 130

fitter() (*d3rlpy.algos.RandomPolicy* method), 160
 fitter() (*d3rlpy.algos.SAC* method), 44
 fitter() (*d3rlpy.algos.TD3* method), 34
 fitter() (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics* method), 327
 fitter() (*d3rlpy.ope.DiscreteFQE* method), 306
 fitter() (*d3rlpy.ope.FQE* method), 296
 FQE (class in *d3rlpy.ope*), 291
 FQFQFunctionFactory (class in *d3rlpy.models.q_functions*), 249
 from_json() (*d3rlpy.algos.AWAC* class method), 109
 from_json() (*d3rlpy.algos.AWR* class method), 99
 from_json() (*d3rlpy.algos.BC* class method), 14
 from_json() (*d3rlpy.algos.BCQ* class method), 56
 from_json() (*d3rlpy.algos.BEAR* class method), 67
 from_json() (*d3rlpy.algos.COMBO* class method), 152
 from_json() (*d3rlpy.algos.CQL* class method), 89
 from_json() (*d3rlpy.algos.CRR* class method), 78
 from_json() (*d3rlpy.algos.DDPG* class method), 24
 from_json() (*d3rlpy.algos.DiscreteAWR* class method), 232
 from_json() (*d3rlpy.algos.DiscreteBC* class method), 171
 from_json() (*d3rlpy.algos.DiscreteBCQ* class method), 212
 from_json() (*d3rlpy.algos.DiscreteCQL* class method), 222
 from_json() (*d3rlpy.algos.DiscreteRandomPolicy* class method), 240
 from_json() (*d3rlpy.algos.DiscreteSAC* class method), 201
 from_json() (*d3rlpy.algos.DoubleDQN* class method), 191
 from_json() (*d3rlpy.algos.DQN* class method), 181
 from_json() (*d3rlpy.algos.MOPO* class method), 141
 from_json() (*d3rlpy.algos.PLAS* class method), 120
 from_json() (*d3rlpy.algos.PLASWithPerturbation* class method), 130
 from_json() (*d3rlpy.algos.RandomPolicy* class method), 161
 from_json() (*d3rlpy.algos.SAC* class method), 45
 from_json() (*d3rlpy.algos.TD3* class method), 35
 from_json() (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics* class method), 327
 from_json() (*d3rlpy.ope.DiscreteFQE* class method), 307
 from_json() (*d3rlpy.ope.FQE* class method), 297
 gamma (*d3rlpy.algos.BEAR* attribute), 71
 gamma (*d3rlpy.algos.COMBO* attribute), 156
 gamma (*d3rlpy.algos.CQL* attribute), 93
 gamma (*d3rlpy.algos.CRR* attribute), 82
 gamma (*d3rlpy.algos.DDPG* attribute), 28
 gamma (*d3rlpy.algos.DiscreteAWR* attribute), 235
 gamma (*d3rlpy.algos.DiscreteBC* attribute), 175
 gamma (*d3rlpy.algos.DiscreteBCQ* attribute), 215
 gamma (*d3rlpy.algos.DiscreteCQL* attribute), 225
 gamma (*d3rlpy.algos.DiscreteRandomPolicy* attribute), 244
 gamma (*d3rlpy.algos.DiscreteSAC* attribute), 205
 gamma (*d3rlpy.algos.DoubleDQN* attribute), 195
 gamma (*d3rlpy.algos.DQN* attribute), 185
 gamma (*d3rlpy.algos.MOPO* attribute), 145
 gamma (*d3rlpy.algos.PLAS* attribute), 123
 gamma (*d3rlpy.algos.PLASWithPerturbation* attribute), 134
 gamma (*d3rlpy.algos.RandomPolicy* attribute), 165
 gamma (*d3rlpy.algos.SAC* attribute), 49
 gamma (*d3rlpy.algos.TD3* attribute), 38
 gamma (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics* attribute), 330
 gamma (*d3rlpy.ope.DiscreteFQE* attribute), 310
 gamma (*d3rlpy.ope.FQE* attribute), 300
 generate_new_data() (*d3rlpy.algos.AWAC* method), 110
 generate_new_data() (*d3rlpy.algos.AWR* method), 100
 generate_new_data() (*d3rlpy.algos.BC* method), 15
 generate_new_data() (*d3rlpy.algos.BCQ* method), 56
 generate_new_data() (*d3rlpy.algos.BEAR* method), 68
 generate_new_data() (*d3rlpy.algos.COMBO* method), 152
 generate_new_data() (*d3rlpy.algos.CQL* method), 89
 generate_new_data() (*d3rlpy.algos.CRR* method), 78
 generate_new_data() (*d3rlpy.algos.DDPG* method), 25
 generate_new_data() (*d3rlpy.algos.DiscreteAWR* method), 232
 generate_new_data() (*d3rlpy.algos.DiscreteBC* method), 172
 generate_new_data() (*d3rlpy.algos.DiscreteBCQ* method), 212
 generate_new_data() (*d3rlpy.algos.DiscreteCQL* method), 222
 generate_new_data() (*d3rlpy.algos.DiscreteRandomPolicy* method), 241

G

gamma (*d3rlpy.algos.AWAC* attribute), 113
 gamma (*d3rlpy.algos.AWR* attribute), 103
 gamma (*d3rlpy.algos.BC* attribute), 18
 gamma (*d3rlpy.algos.BCQ* attribute), 60

`generate_new_data()` (*d3rlpy.algos.DiscreteSAC method*), 202
`generate_new_data()` (*d3rlpy.algos.DoubleDQN method*), 191
`generate_new_data()` (*d3rlpy.algos.DQN method*), 181
`generate_new_data()` (*d3rlpy.algos.MOPO method*), 142
`generate_new_data()` (*d3rlpy.algos.PLAS method*), 120
`generate_new_data()` (*d3rlpy.algos.PLASWithPerturbation method*), 131
`generate_new_data()` (*d3rlpy.algos.RandomPolicy method*), 162
`generate_new_data()` (*d3rlpy.algos.SAC method*), 45
`generate_new_data()` (*d3rlpy.algos.TD3 method*), 35
`generate_new_data()` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 328
`generate_new_data()` (*d3rlpy.ope.DiscreteFQE method*), 307
`generate_new_data()` (*d3rlpy.ope.FQE method*), 297
`get_action_size()` (*d3rlpy.dataset.Episode method*), 257
`get_action_size()` (*d3rlpy.dataset.MDPDataset method*), 254
`get_action_size()` (*d3rlpy.dataset.Transition method*), 258
`get_action_type()` (*d3rlpy.algos.AWAC method*), 110
`get_action_type()` (*d3rlpy.algos.AWR method*), 100
`get_action_type()` (*d3rlpy.algos.BC method*), 15
`get_action_type()` (*d3rlpy.algos.BCQ method*), 57
`get_action_type()` (*d3rlpy.algos.BEAR method*), 68
`get_action_type()` (*d3rlpy.algos.COMBO method*), 153
`get_action_type()` (*d3rlpy.algos.CQL method*), 90
`get_action_type()` (*d3rlpy.algos.CRR method*), 79
`get_action_type()` (*d3rlpy.algos.DDPG method*), 25
`get_action_type()` (*d3rlpy.algos.DiscreteAWR method*), 232
`get_action_type()` (*d3rlpy.algos.DiscreteBC method*), 172
`get_action_type()` (*d3rlpy.algos.DiscreteBCQ method*), 212
`get_action_type()` (*d3rlpy.algos.DiscreteCQL method*), 222
`get_action_type()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 241
`get_action_type()` (*d3rlpy.algos.DiscreteSAC method*), 202
`get_action_type()` (*d3rlpy.algos.DoubleDQN method*), 192
`get_action_type()` (*d3rlpy.algos.DQN method*), 182
`get_action_type()` (*d3rlpy.algos.MOPO method*), 142
`get_action_type()` (*d3rlpy.algos.PLAS method*), 120
`get_action_type()` (*d3rlpy.algos.PLASWithPerturbation method*), 131
`get_action_type()` (*d3rlpy.algos.RandomPolicy method*), 162
`get_action_type()` (*d3rlpy.algos.SAC method*), 46
`get_action_type()` (*d3rlpy.algos.TD3 method*), 35
`get_action_type()` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 328
`get_action_type()` (*d3rlpy.ope.DiscreteFQE method*), 307
`get_action_type()` (*d3rlpy.ope.FQE method*), 297
`get_additional_data()` (*d3rlpy.dataset.TransitionMiniBatch method*), 260
`get_atari()` (in module *d3rlpy.datasets*), 263
`get_cartpole()` (in module *d3rlpy.datasets*), 262
`get_d4rl()` (in module *d3rlpy.datasets*), 264
`get_dataset()` (in module *d3rlpy.datasets*), 264
`get_observation_shape()` (*d3rlpy.dataset.Episode method*), 257
`get_observation_shape()` (*d3rlpy.dataset.MDPDataset method*), 254
`get_observation_shape()` (*d3rlpy.dataset.Transition method*), 258
`get_params()` (*d3rlpy.algos.AWAC method*), 110
`get_params()` (*d3rlpy.algos.AWR method*), 100
`get_params()` (*d3rlpy.algos.BC method*), 15
`get_params()` (*d3rlpy.algos.BCQ method*), 57
`get_params()` (*d3rlpy.algos.BEAR method*), 68
`get_params()` (*d3rlpy.algos.COMBO method*), 153
`get_params()` (*d3rlpy.algos.CQL method*), 90
`get_params()` (*d3rlpy.algos.CRR method*), 79
`get_params()` (*d3rlpy.algos.DDPG method*), 25
`get_params()` (*d3rlpy.algos.DiscreteAWR method*), 232
`get_params()` (*d3rlpy.algos.DiscreteBC method*), 172
`get_params()` (*d3rlpy.algos.DiscreteBCQ method*), 212
`get_params()` (*d3rlpy.algos.DiscreteCQL method*), 222

[get_params\(\)](#) ([d3rlpy.algos.DiscreteCQL method](#)), 222
[get_params\(\)](#) ([d3rlpy.algos.DiscreteRandomPolicy method](#)), 241
[get_params\(\)](#) ([d3rlpy.algos.DiscreteSAC method](#)), 202
[get_params\(\)](#) ([d3rlpy.algos.DoubleDQN method](#)), 192
[get_params\(\)](#) ([d3rlpy.algos.DQN method](#)), 182
[get_params\(\)](#) ([d3rlpy.algos.MOPO method](#)), 142
[get_params\(\)](#) ([d3rlpy.algos.PLAS method](#)), 120
[get_params\(\)](#) ([d3rlpy.algos.PLASWithPerturbation method](#)), 131
[get_params\(\)](#) ([d3rlpy.algos.RandomPolicy method](#)), 162
[get_params\(\)](#) ([d3rlpy.algos.SAC method](#)), 46
[get_params\(\)](#) ([d3rlpy.algos.TD3 method](#)), 35
[get_params\(\)](#) ([d3rlpy.dynamics.ProbabilisticEnsembleDynamics method](#)), 328
[get_params\(\)](#) ([d3rlpy.models.encoders.DefaultEncoderFactory method](#)), 279
[get_params\(\)](#) ([d3rlpy.models.encoders.DenseEncoderFactory method](#)), 283
[get_params\(\)](#) ([d3rlpy.models.encoders.PixelEncoderFactory method](#)), 280
[get_params\(\)](#) ([d3rlpy.models.encoders.VectorEncoderFactory method](#)), 281
[get_params\(\)](#) ([d3rlpy.models.optimizers.AdamFactory method](#)), 275
[get_params\(\)](#) ([d3rlpy.models.optimizers.OptimizerFactory method](#)), 274
[get_params\(\)](#) ([d3rlpy.models.optimizers.RMSpropFactory method](#)), 276
[get_params\(\)](#) ([d3rlpy.models.optimizers.SGDFactory method](#)), 274
[get_params\(\)](#) ([d3rlpy.models.q_functions.FQFQFunctionFactory method](#)), 250
[get_params\(\)](#) ([d3rlpy.models.q_functions.IQNQFunctionFactory method](#)), 249
[get_params\(\)](#) ([d3rlpy.models.q_functions.MeanQFunctionFactory method](#)), 246
[get_params\(\)](#) ([d3rlpy.models.q_functions.QRQFunctionFactory method](#)), 247
[get_params\(\)](#) ([d3rlpy.ope.DiscreteFQE method](#)), 307
[get_params\(\)](#) ([d3rlpy.ope.FQE method](#)), 297
[get_params\(\)](#) ([d3rlpy.preprocessing.MinMaxActionScaler method](#)), 272
[get_params\(\)](#) ([d3rlpy.preprocessing.MinMaxScaler method](#)), 268
[get_params\(\)](#) ([d3rlpy.preprocessing.PixelScaler method](#)), 266
[get_params\(\)](#) ([d3rlpy.preprocessing.StandardScaler method](#)), 270
[get_pendulum\(\)](#) (in module [d3rlpy.datasets](#)), 263
[get_pybullet\(\)](#) (in module [d3rlpy.datasets](#)), 263
[get_type\(\)](#) ([d3rlpy.models.encoders.DefaultEncoderFactory method](#)), 279
[get_type\(\)](#) ([d3rlpy.models.encoders.DenseEncoderFactory method](#)), 283
[get_type\(\)](#) ([d3rlpy.models.encoders.PixelEncoderFactory method](#)), 280
[get_type\(\)](#) ([d3rlpy.models.encoders.VectorEncoderFactory method](#)), 282
[get_type\(\)](#) ([d3rlpy.models.q_functions.FQFQFunctionFactory method](#)), 250
[get_type\(\)](#) ([d3rlpy.models.q_functions.IQNQFunctionFactory method](#)), 249
[get_type\(\)](#) ([d3rlpy.models.q_functions.MeanQFunctionFactory method](#)), 246
[get_type\(\)](#) ([d3rlpy.models.q_functions.QRQFunctionFactory method](#)), 248
[get_type\(\)](#) ([d3rlpy.preprocessing.MinMaxActionScaler method](#)), 272
[get_type\(\)](#) ([d3rlpy.preprocessing.MinMaxScaler method](#)), 268
[get_type\(\)](#) ([d3rlpy.preprocessing.PixelScaler method](#)), 266
[get_type\(\)](#) ([d3rlpy.preprocessing.StandardScaler method](#)), 270
[impl](#) ([d3rlpy.algos.AWAC attribute](#)), 113
[impl](#) ([d3rlpy.algos.AWR attribute](#)), 103
[impl](#) ([d3rlpy.algos.BC attribute](#)), 18
[impl](#) ([d3rlpy.algos.BCQ attribute](#)), 60
[impl](#) ([d3rlpy.algos.BEAR attribute](#)), 71
[impl](#) ([d3rlpy.algos.COMBO attribute](#)), 156
[impl](#) ([d3rlpy.algos.CQL attribute](#)), 93
[impl](#) ([d3rlpy.algos.CRR attribute](#)), 82
[impl](#) ([d3rlpy.algos.DDPG attribute](#)), 28
[impl](#) ([d3rlpy.algos.DiscreteAWR attribute](#)), 235
[impl](#) ([d3rlpy.algos.DiscreteBC attribute](#)), 175
[impl](#) ([d3rlpy.algos.DiscreteBCQ attribute](#)), 215
[impl](#) ([d3rlpy.algos.DiscreteCQL attribute](#)), 225
[impl](#) ([d3rlpy.algos.DiscreteRandomPolicy attribute](#)), 244
[impl](#) ([d3rlpy.algos.DiscreteSAC attribute](#)), 205
[impl](#) ([d3rlpy.algos.DoubleDQN attribute](#)), 195
[impl](#) ([d3rlpy.algos.DQN attribute](#)), 185
[impl](#) ([d3rlpy.algos.MOPO attribute](#)), 145
[impl](#) ([d3rlpy.algos.PLAS attribute](#)), 124
[impl](#) ([d3rlpy.algos.PLASWithPerturbation attribute](#)), 134
[impl](#) ([d3rlpy.algos.RandomPolicy attribute](#)), 165
[impl](#) ([d3rlpy.algos.SAC attribute](#)), 49
[impl](#) ([d3rlpy.algos.TD3 attribute](#)), 38
[impl](#) ([d3rlpy.dynamics.ProbabilisticEnsembleDynamics attribute](#)), 330

impl (*d3rlpy.ope.DiscreteFQE attribute*), 310
 impl (*d3rlpy.ope.FQE attribute*), 300
 initial_state_value_estimation_scorer()
 (*in module d3rlpy.metrics.scorer*), 286
 IQNQFunctionFactory (class *in*
 d3rlpy.models.q_functions), 248
 is_action_discrete()
 (*d3rlpy.dataset.MDPDataset method*), 254

L

LinearDecayEpsilonGreedy (class *in*
 d3rlpy.online.explorers), 319
 load() (*d3rlpy.dataset.MDPDataset class method*), 254
 load_model() (*d3rlpy.algos.AWAC method*), 110
 load_model() (*d3rlpy.algos.AWR method*), 100
 load_model() (*d3rlpy.algos.BC method*), 16
 load_model() (*d3rlpy.algos.BCQ method*), 57
 load_model() (*d3rlpy.algos.BEAR method*), 68
 load_model() (*d3rlpy.algos.COMBO method*), 153
 load_model() (*d3rlpy.algos.CQL method*), 90
 load_model() (*d3rlpy.algos.CRR method*), 79
 load_model() (*d3rlpy.algos.DDPG method*), 25
 load_model() (*d3rlpy.algos.DiscreteAWR method*),
 233
 load_model() (*d3rlpy.algos.DiscreteBC method*),
 172
 load_model() (*d3rlpy.algos.DiscreteBCQ method*),
 213
 load_model() (*d3rlpy.algos.DiscreteCQL method*),
 223
 load_model() (*d3rlpy.algos.DiscreteRandomPolicy*
 method), 241
 load_model() (*d3rlpy.algos.DiscreteSAC method*),
 202
 load_model() (*d3rlpy.algos.DoubleDQN method*),
 192
 load_model() (*d3rlpy.algos.DQN method*), 182
 load_model() (*d3rlpy.algos.MOPO method*), 142
 load_model() (*d3rlpy.algos.PLAS method*), 121
 load_model() (*d3rlpy.algos.PLASWithPerturbation*
 method), 132
 load_model() (*d3rlpy.algos.RandomPolicy method*),
 162
 load_model() (*d3rlpy.algos.SAC method*), 46
 load_model() (*d3rlpy.algos.TD3 method*), 36
 load_model() (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics*
 method), 329
 load_model() (*d3rlpy.ope.DiscreteFQE method*), 308
 load_model() (*d3rlpy.ope.FQE method*), 298

M

mask (*d3rlpy.dataset.Transition attribute*), 259
 masks (*d3rlpy.dataset.TransitionMiniBatch attribute*),
 261

MDPDataset (class *in d3rlpy.dataset*), 252
 MeanQFunctionFactory (class *in*
 d3rlpy.models.q_functions), 246
 MinMaxActionScaler (class *in*
 d3rlpy.preprocessing), 271
 MinMaxScaler (class *in d3rlpy.preprocessing*), 267
 module
 d3rlpy, 9
 d3rlpy.algos, 9
 d3rlpy.dataset, 251
 d3rlpy.datasets, 262
 d3rlpy.dynamics, 323
 d3rlpy.metrics, 283
 d3rlpy.models.encoders, 276
 d3rlpy.models.optimizers, 273
 d3rlpy.models.q_functions, 245
 d3rlpy.online, 316
 d3rlpy.ope, 291
 d3rlpy.preprocessing, 265
 MOPO (class *in d3rlpy.algos*), 135

N

n_frames (*d3rlpy.algos.AWAC attribute*), 113
 n_frames (*d3rlpy.algos.AWR attribute*), 103
 n_frames (*d3rlpy.algos.BC attribute*), 18
 n_frames (*d3rlpy.algos.BCQ attribute*), 60
 n_frames (*d3rlpy.algos.BEAR attribute*), 71
 n_frames (*d3rlpy.algos.COMBO attribute*), 156
 n_frames (*d3rlpy.algos.CQL attribute*), 93
 n_frames (*d3rlpy.algos.CRR attribute*), 82
 n_frames (*d3rlpy.algos.DDPG attribute*), 28
 n_frames (*d3rlpy.algos.DiscreteAWR attribute*), 235
 n_frames (*d3rlpy.algos.DiscreteBC attribute*), 175
 n_frames (*d3rlpy.algos.DiscreteBCQ attribute*), 215
 n_frames (*d3rlpy.algos.DiscreteCQL attribute*), 226
 n_frames (*d3rlpy.algos.DiscreteRandomPolicy at-*
 tribute), 244
 n_frames (*d3rlpy.algos.DiscreteSAC attribute*), 205
 n_frames (*d3rlpy.algos.DoubleDQN attribute*), 195
 n_frames (*d3rlpy.algos.DQN attribute*), 185
 n_frames (*d3rlpy.algos.MOPO attribute*), 145
 n_frames (*d3rlpy.algos.PLAS attribute*), 124
 n_frames (*d3rlpy.algos.PLASWithPerturbation at-*
 tribute), 134
 n_frames (*d3rlpy.algos.RandomPolicy attribute*), 165
 n_frames (*d3rlpy.algos.SAC attribute*), 49
 n_frames (*d3rlpy.algos.TD3 attribute*), 39
 n_frames (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics*
 attribute), 330
 n_frames (*d3rlpy.ope.DiscreteFQE attribute*), 310
 n_frames (*d3rlpy.ope.FQE attribute*), 301
 n_greedy_quantiles
 (*d3rlpy.models.q_functions.IQNQFunctionFactory*
 attribute), 249

`n_quantiles (d3rlpy.models.q_functions.FQFQFunctionFactory attribute)`, 250
`n_quantiles (d3rlpy.models.q_functions.IQNQFunctionFactory attribute)`, 249
`n_quantiles (d3rlpy.models.q_functions.QRQFunctionFactory attribute)`, 248
`n_steps (d3rlpy.algos.AWAC attribute)`, 113
`n_steps (d3rlpy.algos.AWR attribute)`, 103
`n_steps (d3rlpy.algos.BC attribute)`, 18
`n_steps (d3rlpy.algos.BCQ attribute)`, 60
`n_steps (d3rlpy.algos.BEAR attribute)`, 71
`n_steps (d3rlpy.algos.COMBO attribute)`, 156
`n_steps (d3rlpy.algos.CQL attribute)`, 93
`n_steps (d3rlpy.algos.CRR attribute)`, 82
`n_steps (d3rlpy.algos.DDPG attribute)`, 28
`n_steps (d3rlpy.algos.DiscreteAWR attribute)`, 235
`n_steps (d3rlpy.algos.DiscreteBC attribute)`, 175
`n_steps (d3rlpy.algos.DiscreteBCQ attribute)`, 216
`n_steps (d3rlpy.algos.DiscreteCQL attribute)`, 226
`n_steps (d3rlpy.algos.DiscreteRandomPolicy attribute)`, 244
`n_steps (d3rlpy.algos.DiscreteSAC attribute)`, 205
`n_steps (d3rlpy.algos.DoubleDQN attribute)`, 195
`n_steps (d3rlpy.algos.DQN attribute)`, 185
`n_steps (d3rlpy.algos.MOPO attribute)`, 145
`n_steps (d3rlpy.algos.PLAS attribute)`, 124
`n_steps (d3rlpy.algos.PLASWithPerturbation attribute)`, 135
`n_steps (d3rlpy.algos.RandomPolicy attribute)`, 165
`n_steps (d3rlpy.algos.SAC attribute)`, 49
`n_steps (d3rlpy.algos.TD3 attribute)`, 39
`n_steps (d3rlpy.dataset.TransitionMiniBatch attribute)`, 261
`n_steps (d3rlpy.dynamics.ProbabilisticEnsembleDynamics attribute)`, 330
`n_steps (d3rlpy.ope.DiscreteFQE attribute)`, 311
`n_steps (d3rlpy.ope.FQE attribute)`, 301
`next_action (d3rlpy.dataset.Transition attribute)`, 259
`next_actions (d3rlpy.dataset.TransitionMiniBatch attribute)`, 261
`next_observation (d3rlpy.dataset.Transition attribute)`, 259
`next_observations (d3rlpy.dataset.TransitionMiniBatch attribute)`, 261
`next_reward (d3rlpy.dataset.Transition attribute)`, 259
`next_rewards (d3rlpy.dataset.TransitionMiniBatch attribute)`, 261
`next_transition (d3rlpy.dataset.Transition attribute)`, 259
`NormalNoise (class in d3rlpy.online.explorers)`, 320
`observation (d3rlpy.dataset.Transition attribute)`,
`observation_shape (d3rlpy.algos.AWAC attribute)`,
`observation_shape (d3rlpy.algos.AWR attribute)`, 103
`observation_shape (d3rlpy.algos.BC attribute)`, 18
`observation_shape (d3rlpy.algos.BCQ attribute)`, 60
`observation_shape (d3rlpy.algos.BEAR attribute)`, 71
`observation_shape (d3rlpy.algos.COMBO attribute)`, 156
`observation_shape (d3rlpy.algos.CQL attribute)`, 93
`observation_shape (d3rlpy.algos.CRR attribute)`, 82
`observation_shape (d3rlpy.algos.DDPG attribute)`, 29
`observation_shape (d3rlpy.algos.DiscreteAWR attribute)`, 235
`observation_shape (d3rlpy.algos.DiscreteBC attribute)`, 175
`observation_shape (d3rlpy.algos.DiscreteBCQ attribute)`, 216
`observation_shape (d3rlpy.algos.DiscreteCQL attribute)`, 226
`observation_shape (d3rlpy.algos.DiscreteRandomPolicy attribute)`, 245
`observation_shape (d3rlpy.algos.DiscreteSAC attribute)`, 206
`observation_shape (d3rlpy.algos.DoubleDQN attribute)`, 195
`observation_shape (d3rlpy.algos.DQN attribute)`, 185
`observation_shape (d3rlpy.algos.MOPO attribute)`, 146
`observation_shape (d3rlpy.algos.PLAS attribute)`, 124
`observation_shape (d3rlpy.algos.PLASWithPerturbation attribute)`, 135
`observation_shape (d3rlpy.algos.RandomPolicy attribute)`, 166
`observation_shape (d3rlpy.algos.SAC attribute)`, 49
`observation_shape (d3rlpy.algos.TD3 attribute)`, 39
`observation_shape (d3rlpy.dynamics.ProbabilisticEnsembleDynamics attribute)`, 330
`observation_shape (d3rlpy.ope.DiscreteFQE attribute)`,

tribute), 311
 observation_shape (*d3rlpy.ope.FQE attribute*), 301
 observations (*d3rlpy.dataset.Episode attribute*), 257
 observations (*d3rlpy.dataset.MDPDataset attribute*), 255
 observations (*d3rlpy.dataset.TransitionMiniBatch attribute*), 261
 OptimizerFactory (class in *d3rlpy.models.optimizers*), 273

P

PixelEncoderFactory (class in *d3rlpy.models.encoders*), 280
 PixelScaler (class in *d3rlpy.preprocessing*), 266
 PLAS (class in *d3rlpy.algos*), 114
 PLASWithPerturbation (class in *d3rlpy.algos*), 125
 predict() (*d3rlpy.algos.AWAC method*), 111
 predict() (*d3rlpy.algos.AWR method*), 101
 predict() (*d3rlpy.algos.BC method*), 16
 predict() (*d3rlpy.algos.BCQ method*), 57
 predict() (*d3rlpy.algos.BEAR method*), 68
 predict() (*d3rlpy.algos.COMBO method*), 153
 predict() (*d3rlpy.algos.CQL method*), 90
 predict() (*d3rlpy.algos.CRR method*), 79
 predict() (*d3rlpy.algos.DDPG method*), 26
 predict() (*d3rlpy.algos.DiscreteAWR method*), 233
 predict() (*d3rlpy.algos.DiscreteBC method*), 173
 predict() (*d3rlpy.algos.DiscreteBCQ method*), 213
 predict() (*d3rlpy.algos.DiscreteCQL method*), 223
 predict() (*d3rlpy.algos.DiscreteRandomPolicy method*), 242
 predict() (*d3rlpy.algos.DiscreteSAC method*), 203
 predict() (*d3rlpy.algos.DoubleDQN method*), 192
 predict() (*d3rlpy.algos.DQN method*), 182
 predict() (*d3rlpy.algos.MOPO method*), 143
 predict() (*d3rlpy.algos.PLAS method*), 121
 predict() (*d3rlpy.algos.PLASWithPerturbation method*), 132
 predict() (*d3rlpy.algos.RandomPolicy method*), 163
 predict() (*d3rlpy.algos.SAC method*), 46
 predict() (*d3rlpy.algos.TD3 method*), 36
 predict() (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 329
 predict() (*d3rlpy.ope.DiscreteFQE method*), 308
 predict() (*d3rlpy.ope.FQE method*), 298
 predict_value() (*d3rlpy.algos.AWAC method*), 111
 predict_value() (*d3rlpy.algos.AWR method*), 101
 predict_value() (*d3rlpy.algos.BC method*), 16
 predict_value() (*d3rlpy.algos.BCQ method*), 58
 predict_value() (*d3rlpy.algos.BEAR method*), 69
 predict_value() (*d3rlpy.algos.COMBO method*), 153
 predict_value() (*d3rlpy.algos.CQL method*), 90

predict_value() (*d3rlpy.algos.CRR method*), 79
 predict_value() (*d3rlpy.algos.DDPG method*), 26
 predict_value() (*d3rlpy.algos.DiscreteAWR method*), 233
 predict_value() (*d3rlpy.algos.DiscreteBC method*), 173
 predict_value() (*d3rlpy.algos.DiscreteBCQ method*), 213
 predict_value() (*d3rlpy.algos.DiscreteCQL method*), 223
 predict_value() (*d3rlpy.algos.DiscreteRandomPolicy method*), 242
 predict_value() (*d3rlpy.algos.DiscreteSAC method*), 203
 predict_value() (*d3rlpy.algos.DoubleDQN method*), 192
 predict_value() (*d3rlpy.algos.DQN method*), 182
 predict_value() (*d3rlpy.algos.MOPO method*), 143
 predict_value() (*d3rlpy.algos.PLAS method*), 121
 predict_value() (*d3rlpy.algos.PLASWithPerturbation method*), 132
 predict_value() (*d3rlpy.algos.RandomPolicy method*), 163
 predict_value() (*d3rlpy.algos.SAC method*), 47
 predict_value() (*d3rlpy.algos.TD3 method*), 36
 predict_value() (*d3rlpy.ope.DiscreteFQE method*), 308
 predict_value() (*d3rlpy.ope.FQE method*), 298
 prev_transition (*d3rlpy.dataset.Transition attribute*), 259
 ProbabilisticEnsembleDynamics (class in *d3rlpy.dynamics*), 324

Q

QRQFunctionFactory (class in *d3rlpy.models.q_functions*), 247

R

RandomPolicy (class in *d3rlpy.algos*), 157
 ReplayBuffer (class in *d3rlpy.online.buffers*), 317
 reverse_transform() (*d3rlpy.preprocessing.MinMaxActionScaler method*), 272
 reverse_transform() (*d3rlpy.preprocessing.MinMaxScaler method*), 268
 reverse_transform() (*d3rlpy.preprocessing.PixelScaler method*), 267
 reverse_transform() (*d3rlpy.preprocessing.StandardScaler method*), 270

reverse_transform_numpy() (d3rlpy.preprocessing.MinMaxActionScaler method), 272
 reward (d3rlpy.dataset.Transition attribute), 259
 rewards (d3rlpy.dataset.Episode attribute), 257
 rewards (d3rlpy.dataset.MDPDataset attribute), 255
 rewards (d3rlpy.dataset.TransitionMiniBatch attribute), 262
 RMSpropFactory (class in d3rlpy.models.optimizers), 276

S

SAC (class in d3rlpy.algos), 39
 sample() (d3rlpy.online.buffers.BatchReplayBuffer method), 322
 sample() (d3rlpy.online.buffers.ReplayBuffer method), 318
 sample() (d3rlpy.online.explorers.ConstantEpsilonGreedy method), 319
 sample() (d3rlpy.online.explorers.LinearDecayEpsilonGreedy method), 320
 sample() (d3rlpy.online.explorers.NormalNoise method), 320
 sample_action() (d3rlpy.algos.AWAC method), 111
 sample_action() (d3rlpy.algos.AWR method), 101
 sample_action() (d3rlpy.algos.BC method), 16
 sample_action() (d3rlpy.algos.BCQ method), 58
 sample_action() (d3rlpy.algos.BEAR method), 69
 sample_action() (d3rlpy.algos.COMBO method), 154
 sample_action() (d3rlpy.algos.CQL method), 91
 sample_action() (d3rlpy.algos.CRR method), 80
 sample_action() (d3rlpy.algos.DDPG method), 26
 sample_action() (d3rlpy.algos.DiscreteAWR method), 233
 sample_action() (d3rlpy.algos.DiscreteBC method), 173
 sample_action() (d3rlpy.algos.DiscreteBCQ method), 214
 sample_action() (d3rlpy.algos.DiscreteCQL method), 224
 sample_action() (d3rlpy.algos.DiscreteRandomPolicy method), 242
 sample_action() (d3rlpy.algos.DiscreteSAC method), 203
 sample_action() (d3rlpy.algos.DoubleDQN method), 193
 sample_action() (d3rlpy.algos.DQN method), 183
 sample_action() (d3rlpy.algos.MOPO method), 143
 sample_action() (d3rlpy.algos.PLAS method), 122
 sample_action() (d3rlpy.algos.PLASWithPerturbation method), 133
 sample_action() (d3rlpy.algos.RandomPolicy method), 163
 sample_action() (d3rlpy.algos.SAC method), 47
 sample_action() (d3rlpy.algos.TD3 method), 37
 sample_action() (d3rlpy.ope.DiscreteFQE method), 309
 sample_action() (d3rlpy.ope.FQE method), 299
 save_model() (d3rlpy.algos.AWAC method), 112
 save_model() (d3rlpy.algos.AWR method), 101
 save_model() (d3rlpy.algos.BC method), 16
 save_model() (d3rlpy.algos.BCQ method), 58
 save_model() (d3rlpy.algos.BEAR method), 69
 save_model() (d3rlpy.algos.COMBO method), 154
 save_model() (d3rlpy.algos.CQL method), 91
 save_model() (d3rlpy.algos.CRR method), 80
 save_model() (d3rlpy.algos.DDPG method), 27
 save_model() (d3rlpy.algos.DiscreteAWR method), 234
 save_model() (d3rlpy.algos.DiscreteBC method), 173
 save_model() (d3rlpy.algos.DiscreteBCQ method), 214
 save_model() (d3rlpy.algos.DiscreteCQL method), 224
 save_model() (d3rlpy.algos.DiscreteRandomPolicy method), 243
 save_model() (d3rlpy.algos.DiscreteSAC method), 204
 save_model() (d3rlpy.algos.DoubleDQN method), 193
 save_model() (d3rlpy.algos.DQN method), 183
 save_model() (d3rlpy.algos.MOPO method), 144
 save_model() (d3rlpy.algos.PLAS method), 122
 save_model() (d3rlpy.algos.PLASWithPerturbation method), 133
 save_model() (d3rlpy.algos.RandomPolicy method), 164
 save_model() (d3rlpy.algos.SAC method), 47
 save_model() (d3rlpy.algos.TD3 method), 37
 save_model() (d3rlpy.dynamics.ProbabilisticEnsembleDynamics method), 329
 save_model() (d3rlpy.ope.DiscreteFQE method), 309
 save_model() (d3rlpy.ope.FQE method), 299
 save_params() (d3rlpy.algos.AWAC method), 112
 save_params() (d3rlpy.algos.AWR method), 101
 save_params() (d3rlpy.algos.BC method), 16
 save_params() (d3rlpy.algos.BCQ method), 58
 save_params() (d3rlpy.algos.BEAR method), 70
 save_params() (d3rlpy.algos.COMBO method), 154
 save_params() (d3rlpy.algos.CQL method), 91
 save_params() (d3rlpy.algos.CRR method), 80
 save_params() (d3rlpy.algos.DDPG method), 27
 save_params() (d3rlpy.algos.DiscreteAWR method), 234

`save_params()` (*d3rlpy.algos.DiscreteBC method*), 173
`save_params()` (*d3rlpy.algos.DiscreteBCQ method*), 214
`save_params()` (*d3rlpy.algos.DiscreteCQL method*), 224
`save_params()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 243
`save_params()` (*d3rlpy.algos.DiscreteSAC method*), 204
`save_params()` (*d3rlpy.algos.DoubleDQN method*), 193
`save_params()` (*d3rlpy.algos.DQN method*), 183
`save_params()` (*d3rlpy.algos.MOPO method*), 144
`save_params()` (*d3rlpy.algos.PLAS method*), 122
`save_params()` (*d3rlpy.algos.PLASWithPerturbation method*), 133
`save_params()` (*d3rlpy.algos.RandomPolicy method*), 164
`save_params()` (*d3rlpy.algos.SAC method*), 47
`save_params()` (*d3rlpy.algos.TD3 method*), 37
`save_params()` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamic method*), 329
`save_params()` (*d3rlpy.ope.DiscreteFQE method*), 309
`save_params()` (*d3rlpy.ope.FQE method*), 299
`save_policy()` (*d3rlpy.algos.AWAC method*), 112
`save_policy()` (*d3rlpy.algos.AWR method*), 102
`save_policy()` (*d3rlpy.algos.BC method*), 17
`save_policy()` (*d3rlpy.algos.BCQ method*), 58
`save_policy()` (*d3rlpy.algos.BEAR method*), 70
`save_policy()` (*d3rlpy.algos.COMBO method*), 154
`save_policy()` (*d3rlpy.algos.CQL method*), 91
`save_policy()` (*d3rlpy.algos.CRR method*), 80
`save_policy()` (*d3rlpy.algos.DDPG method*), 27
`save_policy()` (*d3rlpy.algos.DiscreteAWR method*), 234
`save_policy()` (*d3rlpy.algos.DiscreteBC method*), 173
`save_policy()` (*d3rlpy.algos.DiscreteBCQ method*), 214
`save_policy()` (*d3rlpy.algos.DiscreteCQL method*), 224
`save_policy()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 243
`save_policy()` (*d3rlpy.algos.DiscreteSAC method*), 204
`save_policy()` (*d3rlpy.algos.DoubleDQN method*), 193
`save_policy()` (*d3rlpy.algos.DQN method*), 183
`save_policy()` (*d3rlpy.algos.MOPO method*), 144
`save_policy()` (*d3rlpy.algos.PLAS method*), 122
`save_policy()` (*d3rlpy.algos.PLASWithPerturbation method*), 133
`save_policy()` (*d3rlpy.algos.RandomPolicy method*), 164
`save_policy()` (*d3rlpy.algos.SAC method*), 47
`save_policy()` (*d3rlpy.algos.TD3 method*), 37
`save_policy()` (*d3rlpy.ope.DiscreteFQE method*), 309
`save_policy()` (*d3rlpy.ope.FQE method*), 299
`scaler` (*d3rlpy.algos.AWAC attribute*), 114
`scaler` (*d3rlpy.algos.AWR attribute*), 103
`scaler` (*d3rlpy.algos.BC attribute*), 18
`scaler` (*d3rlpy.algos.BCQ attribute*), 60
`scaler` (*d3rlpy.algos.BEAR attribute*), 71
`scaler` (*d3rlpy.algos.COMBO attribute*), 156
`scaler` (*d3rlpy.algos.CQL attribute*), 93
`scaler` (*d3rlpy.algos.CRR attribute*), 82
`scaler` (*d3rlpy.algos.DDPG attribute*), 29
`scaler` (*d3rlpy.algos.DiscreteAWR attribute*), 236
`scaler` (*d3rlpy.algos.DiscreteBC attribute*), 175
`scaler` (*d3rlpy.algos.DiscreteBCQ attribute*), 216
`scaler` (*d3rlpy.algos.DiscreteCQL attribute*), 226
`scaler` (*d3rlpy.algos.DiscreteRandomPolicy attribute*), 245
`scaler` (*d3rlpy.algos.DiscreteSAC attribute*), 206
`scaler` (*d3rlpy.algos.DoubleDQN attribute*), 195
`scaler` (*d3rlpy.algos.DQN attribute*), 185
`scaler` (*d3rlpy.algos.MOPO attribute*), 146
`scaler` (*d3rlpy.algos.PLAS attribute*), 124
`scaler` (*d3rlpy.algos.PLASWithPerturbation attribute*), 135
`scaler` (*d3rlpy.algos.RandomPolicy attribute*), 166
`scaler` (*d3rlpy.algos.SAC attribute*), 49
`scaler` (*d3rlpy.algos.TD3 attribute*), 39
`scaler` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics attribute*), 331
`scaler` (*d3rlpy.ope.DiscreteFQE attribute*), 311
`scaler` (*d3rlpy.ope.FQE attribute*), 301
`set_params()` (*d3rlpy.algos.AWAC method*), 112
`set_params()` (*d3rlpy.algos.AWR method*), 102
`set_params()` (*d3rlpy.algos.BC method*), 17
`set_params()` (*d3rlpy.algos.BCQ method*), 59
`set_params()` (*d3rlpy.algos.BEAR method*), 70
`set_params()` (*d3rlpy.algos.COMBO method*), 155
`set_params()` (*d3rlpy.algos.CQL method*), 92
`set_params()` (*d3rlpy.algos.CRR method*), 81
`set_params()` (*d3rlpy.algos.DDPG method*), 27
`set_params()` (*d3rlpy.algos.DiscreteAWR method*), 234
`set_params()` (*d3rlpy.algos.DiscreteBC method*), 174
`set_params()` (*d3rlpy.algos.DiscreteBCQ method*), 214
`set_params()` (*d3rlpy.algos.DiscreteCQL method*), 225

- set_params() (*d3rlpy.algos.DiscreteRandomPolicy method*), 243
 set_params() (*d3rlpy.algos.DiscreteSAC method*), 204
 set_params() (*d3rlpy.algos.DoubleDQN method*), 194
 set_params() (*d3rlpy.algos.DQN method*), 184
 set_params() (*d3rlpy.algos.MOPO method*), 144
 set_params() (*d3rlpy.algos.PLAS method*), 123
 set_params() (*d3rlpy.algos.PLASWithPerturbation method*), 133
 set_params() (*d3rlpy.algos.RandomPolicy method*), 164
 set_params() (*d3rlpy.algos.SAC method*), 48
 set_params() (*d3rlpy.algos.TD3 method*), 38
 set_params() (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 329
 set_params() (*d3rlpy.ope.DiscreteFQE method*), 309
 set_params() (*d3rlpy.ope.FQE method*), 300
 SGDFactory (*class in d3rlpy.models.optimizers*), 274
 share_encoder(*d3rlpy.models.q_functions.FQFQFunctionFactory attribute*), 250
 share_encoder(*d3rlpy.models.q_functions.IQNQFunctionFactory attribute*), 249
 share_encoder(*d3rlpy.models.q_functions.MeanQFunctionFactory attribute*), 247
 share_encoder(*d3rlpy.models.q_functions.QRQFunctionFactory attribute*), 248
 size() (*d3rlpy.dataset.Episode method*), 257
 size() (*d3rlpy.dataset.MDPDataset method*), 255
 size() (*d3rlpy.dataset.TransitionMiniBatch method*), 261
 size() (*d3rlpy.online.buffers.BatchReplayBuffer method*), 323
 size() (*d3rlpy.online.buffers.ReplayBuffer method*), 318
 soft_opc_scorer() (*in module d3rlpy.metrics.scorer*), 286
 StandardScaler (*class in d3rlpy.preprocessing*), 269
- ## T
- TD3 (*class in d3rlpy.algos*), 29
 td_error_scorer() (*in module d3rlpy.metrics.scorer*), 284
 terminal (*d3rlpy.dataset.Episode attribute*), 257
 terminal (*d3rlpy.dataset.Transition attribute*), 259
 terminals (*d3rlpy.dataset.MDPDataset attribute*), 255
 terminals (*d3rlpy.dataset.TransitionMiniBatch attribute*), 262
 to_mdp_dataset() (*d3rlpy.online.buffers.BatchReplayBuffer method*), 323
 to_mdp_dataset() (*d3rlpy.online.buffers.ReplayBuffer method*), 318
- transform() (*d3rlpy.preprocessing.MinMaxActionScaler method*), 272
 transform() (*d3rlpy.preprocessing.MinMaxScaler method*), 268
 transform() (*d3rlpy.preprocessing.PixelScaler method*), 267
 transform() (*d3rlpy.preprocessing.StandardScaler method*), 270
 Transition (*class in d3rlpy.dataset*), 258
 TransitionMiniBatch (*class in d3rlpy.dataset*), 260
 transitions (*d3rlpy.dataset.Episode attribute*), 257
 transitions (*d3rlpy.dataset.TransitionMiniBatch attribute*), 262
 transitions (*d3rlpy.online.buffers.BatchReplayBuffer attribute*), 323
 transitions (*d3rlpy.online.buffers.ReplayBuffer attribute*), 319
 TYPE (*d3rlpy.models.encoders.DefaultEncoderFactory attribute*), 280
 TYPE (*d3rlpy.models.encoders.DenseEncoderFactory attribute*), 283
 TYPE (*d3rlpy.models.encoders.PixelEncoderFactory attribute*), 281
 TYPE (*d3rlpy.models.encoders.VectorEncoderFactory attribute*), 282
 TYPE (*d3rlpy.models.q_functions.FQFQFunctionFactory attribute*), 250
 TYPE (*d3rlpy.models.q_functions.IQNQFunctionFactory attribute*), 249
 TYPE (*d3rlpy.models.q_functions.MeanQFunctionFactory attribute*), 247
 TYPE (*d3rlpy.models.q_functions.QRQFunctionFactory attribute*), 248
 TYPE (*d3rlpy.preprocessing.MinMaxActionScaler attribute*), 273
 TYPE (*d3rlpy.preprocessing.MinMaxScaler attribute*), 269
 TYPE (*d3rlpy.preprocessing.PixelScaler attribute*), 267
 TYPE (*d3rlpy.preprocessing.StandardScaler attribute*), 270
- ## U
- update() (*d3rlpy.algos.AWAC method*), 113
 update() (*d3rlpy.algos.AWR method*), 102
 update() (*d3rlpy.algos.BC method*), 17
 update() (*d3rlpy.algos.BCQ method*), 59
 update() (*d3rlpy.algos.BEAR method*), 70
 update() (*d3rlpy.algos.COMBO method*), 155
 update() (*d3rlpy.algos.CQL method*), 92
 update() (*d3rlpy.algos.CRR method*), 81
 update() (*d3rlpy.algos.DDPG method*), 28
 update() (*d3rlpy.algos.DiscreteAWR method*), 235
 update() (*d3rlpy.algos.DiscreteBC method*), 174

`update()` (*d3rlpy.algos.DiscreteBCQ method*), 215
`update()` (*d3rlpy.algos.DiscreteCQL method*), 225
`update()` (*d3rlpy.algos.DiscreteRandomPolicy method*), 244
`update()` (*d3rlpy.algos.DiscreteSAC method*), 205
`update()` (*d3rlpy.algos.DoubleDQN method*), 194
`update()` (*d3rlpy.algos.DQN method*), 184
`update()` (*d3rlpy.algos.MOPO method*), 145
`update()` (*d3rlpy.algos.PLAS method*), 123
`update()` (*d3rlpy.algos.PLASWithPerturbation method*), 134
`update()` (*d3rlpy.algos.RandomPolicy method*), 165
`update()` (*d3rlpy.algos.SAC method*), 48
`update()` (*d3rlpy.algos.TD3 method*), 38
`update()` (*d3rlpy.dynamics.ProbabilisticEnsembleDynamics method*), 329
`update()` (*d3rlpy.ope.DiscreteFQE method*), 310
`update()` (*d3rlpy.ope.FQE method*), 300

V

`value_estimation_std_scorer()` (in module *d3rlpy.metrics.scorer*), 285
`VectorEncoderFactory` (class in *d3rlpy.models.encoders*), 281