

---

# MushroomRL Documentation

*Release 1.4.0*

**Carlo D'Eramo, Davide Tateo**

**Apr 22, 2020**



<b>1 Reinforcement Learning python library</b>	<b>1</b>
<b>2 Basic run example</b>	<b>3</b>
<b>3 Download and installation</b>	<b>5</b>
3.1 Agent-Environment Interface . . . . .	5
3.2 Actor-Critic . . . . .	10
3.3 Policy search . . . . .	24
3.4 Value-Based . . . . .	33
3.5 Approximators . . . . .	55
3.6 Distributions . . . . .	60
3.7 Environments . . . . .	65
3.8 Features . . . . .	89
3.9 Policy . . . . .	94
3.10 Solvers . . . . .	111
3.11 Utils . . . . .	112
3.12 How to make a simple experiment . . . . .	135
3.13 How to make an advanced experiment . . . . .	136
3.14 How to create a regressor . . . . .	138
3.15 How to make a deep RL experiment . . . . .	140
<b>Python Module Index</b>	<b>147</b>
<b>Index</b>	<b>149</b>



# CHAPTER 1

---

## Reinforcement Learning python library

---

MushroomRL is a Reinforcement Learning (RL) library that aims to be a simple, yet powerful way to make **RL** and **deep RL** experiments. The idea behind Mushroom consists in offering the majority of RL algorithms providing a common interface in order to run them without excessive effort. Moreover, it is designed in such a way that new algorithms and other stuff can generally be added transparently without the need of editing other parts of the code. MushroomRL makes a large use of the environments provided by [OpenAI Gym](#), [DeepMind Control Suite](#) and [MuJoCo](#) libraries, and the [PyTorch](#) library for tensor computation.

With MushroomRL you can:

- solve RL problems simply writing a single small script;
- add custom algorithms and other stuff transparently;
- use all RL environments offered by well-known libraries and build customized environments as well;
- exploit regression models offered by Scikit-Learn or build a customized one with PyTorch;
- run experiments on GPU.



# CHAPTER 2

---

## Basic run example

---

Solve a discrete MDP in few a lines. Firstly, create a **MDP**:

```
from mushroom_rl.environments import GridWorld
mdp = GridWorld(width=3, height=3, goal=(2, 2), start=(0, 0))
```

Then, an epsilon-greedy **policy** with:

```
from mushroom_rl.policy import EpsGreedy
from mushroom_rl.utils.parameters import Parameter

epsilon = Parameter(value=1.)
policy = EpsGreedy(epsilon=epsilon)
```

Eventually, the **agent** is:

```
from mushroom_rl.algorithms.value import QLearning

learning_rate = Parameter(value=.6)
agent = QLearning(policy, mdp.info, learning_rate)
```

Learn:

```
from mushroom_rl.core.core import Core

core = Core(agent, mdp)
core.learn(n_steps=10000, n_steps_per_fit=1)
```

Print final Q-table:

```
import numpy as np

shape = agent.approximator.shape
q = np.zeros(shape)
```

(continues on next page)

(continued from previous page)

```
for i in range(shape[0]):  
    for j in range(shape[1]):  
        state = np.array([i])  
        action = np.array([j])  
        q[i, j] = agent.approximator.predict(state, action)  
print(q)
```

Results in:

```
[[ 6.561   7.29   6.561   7.29 ]  
 [ 7.29    8.1    6.561   8.1   ]  
 [ 8.1     9.     7.29   8.1   ]  
 [ 6.561   8.1    7.29   8.1   ]  
 [ 7.29    9.     7.29   9.    ]  
 [ 8.1     10.    8.1    9.    ]  
 [ 7.29    8.1    8.1    9.    ]  
 [ 8.1     9.     8.1    10.   ]  
 [ 0.      0.     0.     0.    ]]
```

where the Q-values of each action of the MDP are stored for each rows representing a state of the MDP.

# CHAPTER 3

---

## Download and installation

---

MushroomRL can be downloaded from the [GitHub](#) repository. Installation can be done running

```
pip3 install mushroom_rl
```

To compile the documentation:

```
cd mushroom_rl/docs  
make html
```

or to compile the pdf version:

```
cd mushroom_rl/docs  
make latexpdf
```

To launch MushroomRL test suite:

```
pytest
```

## 3.1 Agent-Environment Interface

The three basic interface of mushroom\_rl are the Agent, the Environment and the Core interface.

- The `Agent` is the basic interface for any Reinforcement Learning algorithm.
- The `Environment` is the basic interface for every problem/task that the agent should solve.
- The `Core` is a class used to control the interaction between an agent and an environment.

### 3.1.1 Agent

MushroomRL provides the implementations of several algorithms belonging to all categories of RL:

- value-based;
- policy-search;
- actor-critic.

One can easily implement customized algorithms following the structure of the already available ones, by extending the following interface:

**class** mushroom\_rl.algorithms.agent.**Agent** (*mdp\_info*, *policy*, *features=None*)

Bases: object

This class implements the functions to manage the agent (e.g. move the agent following its policy).

**\_\_init\_\_** (*mdp\_info*, *policy*, *features=None*)

Constructor.

### Parameters

- **mdp\_info** (`MDPInfo`) – information about the MDP;
- **policy** (`Policy`) – the policy followed by the agent;
- **features** (`object`, `None`) – features to extract from the state.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start** ()

Called by the agent when a new episode starts.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**classmethod load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**copy** ()

**Returns** A deepcopy of the agent.

**\_add\_save\_attr** (\*\**attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

### 3.1.2 Environment

MushroomRL provides several implementation of well known benchmarks with both continuous and discrete action spaces.

To implement a new environment, it is mandatory to use the following interface:

```
class mushroom_rl.environments.environment.MDPInfo (observation_space, action_space,  
                                         gamma, horizon)
```

Bases: object

This class is used to store the information of the environment.

```
__init__ (observation_space, action_space, gamma, horizon)  
    Constructor.
```

#### Parameters

- **observation\_space** ([*Box*, *Discrete*]) – the state space;
- **action\_space** ([*Box*, *Discrete*]) – the action space;
- **gamma** (*float*) – the discount factor;
- **horizon** (*int*) – the horizon.

#### **size**

The sum of the number of discrete states and discrete actions. Only works for discrete spaces.

**Type** Returns

#### **shape**

The concatenation of the shape tuple of the state and action spaces.

**Type** Returns

```
class mushroom_rl.environments.environment.Environment (mdp_info)
```

Bases: object

Basic interface used by any mushroom environment.

```
__init__ (mdp_info)  
    Constructor.
```

**Parameters** **mdp\_info** (*MDPInfo*) – an object containing the info of the environment.

#### **seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

#### **reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray*, *None*) – the state to set to the current state.

**Returns** The current state.

#### **step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

### `stop()`

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

### `info`

An object containing the info of the environment.

**Type** Returns

### `static _bound(x, min_value, max_value)`

Method used to bound state and action variables.

#### Parameters

- `x` – the variable to bound;
- `min_value` – the minimum value;
- `max_value` – the maximum value;

**Returns** The bounded variable.

### 3.1.3 Core

```
class mushroom_rl.core.core.Core(agent, mdp, callbacks_episode=None, callback_step=None,
                                 preprocessors=None)
```

Bases: `object`

Implements the functions to run a generic algorithm.

### `__init__(agent, mdp, callbacks_episode=None, callback_step=None, preprocessors=None)`

Constructor.

#### Parameters

- `agent` ([Agent](#)) – the agent moving according to a policy;
- `mdp` ([Environment](#)) – the environment in which the agent moves;
- `callbacks_episode` (`list`) – list of callbacks to execute at the end of each learn iteration;
- `callback_step` ([Callback](#)) – callback to execute after each step;
- `preprocessors` (`list`) – list of state preprocessors to be applied to state variables before feeding them to the agent.

### `learn(n_steps=None, n_episodes=None, n_steps_per_fit=None, n_episodes_per_fit=None, render=False, quiet=False)`

This function moves the agent in the environment and fits the policy using the collected samples. The agent can be moved for a given number of steps or a given number of episodes and, independently from this choice, the policy can be fitted after a given number of steps or a given number of episodes. By default, the environment is reset.

#### Parameters

- `n_steps` (`int, None`) – number of steps to move the agent;
- `n_episodes` (`int, None`) – number of episodes to move the agent;

- **n\_steps\_per\_fit** (*int, None*) – number of steps between each fit of the policy;
- **n\_episodes\_per\_fit** (*int, None*) – number of episodes between each fit of the policy;
- **render** (*bool, False*) – whether to render the environment or not;
- **quiet** (*bool, False*) – whether to show the progress bar or not.

**evaluate** (*initial\_states=None, n\_steps=None, n\_episodes=None, render=False, quiet=False*)

This function moves the agent in the environment using its policy. The agent is moved for a provided number of steps, episodes, or from a set of initial states for the whole episode. By default, the environment is reset.

#### Parameters

- **initial\_states** (*np.ndarray, None*) – the starting states of each episode;
- **n\_steps** (*int, None*) – number of steps to move the agent;
- **n\_episodes** (*int, None*) – number of episodes to move the agent;
- **render** (*bool, False*) – whether to render the environment or not;
- **quiet** (*bool, False*) – whether to show the progress bar or not.

**\_step** (*render*)

Single step.

**Parameters** **render** (*bool*) – whether to render or not.

**Returns** A tuple containing the previous state, the action sampled by the agent, the reward obtained, the reached state, the absorbing flag of the reached state and the last step flag.

**reset** (*initial\_states=None*)

Reset the state of the agent.

**\_preprocess** (*state*)

Method to apply state preprocessors.

**Parameters** **state** (*np.ndarray*) – the state to be preprocessed.

**Returns** The preprocessed state.

## 3.2 Actor-Critic

### 3.2.1 Classical Actor-Critic Methods

```
class mushroom_rl.algorithms.actor_critic.classic_actor_critic.COPDAC_Q(mdp_info,
    pol-
    icy,
    mu,
    al-
    pha_theta,
    al-
    pha_omega,
    al-
    pha_v,
    value_function_features=None,
    pol-
    icy_features=None)
```

Bases: *mushroom\_rl.algorithms.agent.Agent*

Compatible off-policy deterministic actor-critic algorithm. “Deterministic Policy Gradient Algorithms”. Silver D. et al.. 2014.

```
__init__(mdp_info, policy, mu, alpha_theta, alpha_omega, alpha_v, value_function_features=None,
        policy_features=None)
```

Constructor.

#### Parameters

- **mu** (`Regressor`) – regressor that describe the deterministic policy to be learned i.e., the deterministic mapping between state and action.
- **alpha\_theta** (`Parameter`) – learning rate for policy update;
- **alpha\_omega** (`Parameter`) – learning rate for the advantage function;
- **alpha\_v** (`Parameter`) – learning rate for the value function;
- **value\_function\_features** (`Features, None`) – features used by the value function approximator;
- **policy\_features** (`Features, None`) – features used by the policy.

```
fit(dataset)
```

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

```
_add_save_attr(**attr_dict)
```

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

```
_post_load()
```

This method can be overwritten to implement logic that is executed after the loading of the agent.

```
copy()
```

**Returns** A deepcopy of the agent.

**draw\_action(state)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state (np.ndarray)` – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** `path (string)` – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** `path (string)` – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.classic_actor_critic.StochasticAC(mdp_info,
    pol-
    icy,
    al-
    pha_theta,
    al-
    pha_v,
    lambda_par=0.9,
    value_function_featu-
    pol-
    icy_features=None)
```

Bases: `mushroom_rl.algorithms.agent.Agent`

Stochastic Actor critic in the episodic setting as presented in: “Model-Free Reinforcement Learning with Continuous Action in Practice”. Degris T. et al.. 2012.

```
__init__(mdp_info, policy, alpha_theta, alpha_v, lambda_par=0.9, value_function_features=None,
        policy_features=None)
```

Constructor.

**Parameters**

- `alpha_theta` (`Parameter`) – learning rate for policy update;
- `alpha_v` (`Parameter`) – learning rate for the value function;
- `lambda_par` (`float, 0`) – trace decay parameter;
- `value_function_features` (`Features, None`) – features used by the value function approximator;
- `policy_features` (`Features, None`) – features used by the policy.

**episode\_start()**

Called by the agent when a new episode starts.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**\_add\_save\_attr** (*\*\*attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load** ()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**classmethod** **load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.classic_actor_critic.StochasticAC_AVG (mdp_info,
    pol-
    icy,
    al-
    pha_theta,
    al-
    pha_v,
    al-
    pha_r,
    lambda_par=0,
    value_function_
    pol-
    icy_features=N)
```

Bases: mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic.  
stochastic\_ac.StochasticAC

Stochastic Actor critic in the average reward setting as presented in: “Model-Free Reinforcement Learning with Continuous Action in Practice”. Degris T. et al.. 2012.

---

```
__init__(mdp_info, policy, alpha_theta, alpha_v, alpha_r, lambda_par=0.9,
        value_function_features=None, policy_features=None)
    Constructor.

    Parameters alpha_r (Parameter) – learning rate for the reward trace.

__add_save_attr__(**attr_dict)
    Add attributes that should be saved for an agent.

    Parameters attr_dict (dict) – dictionary of attributes mapped to the method that should
        be used to save and load them.

_post_load()
    This method can be overwritten to implement logic that is executed after the loading of the agent.

copy()
    Returns A deepcopy of the agent.

draw_action(state)
    Return the action to execute in the given state. It is the action returned by the policy or the action set by
    the algorithm (e.g. in the case of SARSA).

    Parameters state (np.ndarray) – the state where the agent is.

    Returns The action to be executed.

episode_start()
    Called by the agent when a new episode starts.

fit(dataset)
    Fit step.

    Parameters dataset (list) – the dataset.

classmethod load(path)
    Load and deserialize the agent from the given location on disk.

    Parameters path (string) – Relative or absolute path to the agents save location.

    Returns The loaded agent.

save(path)
    Serialize and save the agent to the given path on disk.

    Parameters path (string) – Relative or absolute path to the agents save location.

stop()
    Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup
    environments internals after a core learn/evaluate to enforce consistency.
```

### 3.2.2 Deep Actor-Critic Methods

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.DeepAC(mdp_info,
        policy, actor_optimizer,
        parameters)
```

Bases: `mushroom_rl.algorithms.agent.Agent`

Base class for algorithms that uses the reparametrization trick, such as SAC, DDPG and TD3.

```
__init__(mdp_info, policy, actor_optimizer, parameters)
    Constructor.
```

**Parameters**

- **actor\_optimizer** (*dict*) – parameters to specify the actor optimizer algorithm;
- **parameters** – policy parameters to be optimized.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.**\_optimize\_actor\_parameters** (*loss*)

Method used to update actor parameters to maximize a given loss.

**Parameters** **loss** (*torch.tensor*) – the loss computed by the algorithm.**\_add\_save\_attr** (*\*\*attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.**\_post\_load** ()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()**Returns** A deepcopy of the agent.**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.**Returns** The action to be executed.**episode\_start** ()

Called by the agent when a new episode starts.

**classmethod** **load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.**Returns** The loaded agent.**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.A2C(mdp_info,  
                                         policy,      ac-  
                                         tor_optimizer,  
                                         critic_params,  
                                         ent_coeff,  
                                         max_grad_norm=None,  
                                         critic_fit_params=None)
```

Bases: mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic.  
deep\_actor\_critic.DeepAC

Advantage Actor Critic algorithm (A2C). Synchronous version of the A3C algorithm. “Asynchronous Methods for Deep Reinforcement Learning”. Mnih V. et. al.. 2016.

**\_\_init\_\_(self, mdp\_info, policy, actor\_optimizer, critic\_params, ent\_coeff, max\_grad\_norm=None, critic\_fit\_params=None)**  
Constructor.

#### Parameters

- **policy** (`TorchPolicy`) – torch policy to be learned by the algorithm;
- **actor\_optimizer** (`dict`) – parameters to specify the actor optimizer algorithm;
- **critic\_params** (`dict`) – parameters of the critic approximator to build;
- **ent\_coeff** (`float, 0`) – coefficient for the entropy penalty;
- **max\_grad\_norm** (`float, None`) – maximum norm for gradient clipping. If `None`, no clipping will be performed, unless specified otherwise in `actor_optimizer`;
- **critic\_fit\_params** (`dict, None`) – parameters of the fitting algorithm of the critic approximator.

**fit(self, dataset)**

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_optimize\_actor\_parameters(loss)**

Method used to update actor parameters to maximize a given loss.

**Parameters** `loss` (`torch.tensor`) – the loss computed by the algorithm.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action(state)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.DDPG(mdp_info,
                                                               policy_class,
                                                               pol-
                                                               icy_params,
                                                               ac-
                                                               tor_params,
                                                               ac-
                                                               tor_optimizer,
                                                               critic_params,
                                                               batch_size,
                                                               ini-
                                                               tial_replay_size,
                                                               max_replay_size,
                                                               tau,      pol-
                                                               icy_delay=1,
                                                               critic_fit_params=None)
```

Bases: `mushroom_rl.algorithms.actor_critic.deep_actor_critic.  
deep_actor_critic.DeepAC`

Deep Deterministic Policy Gradient algorithm. “Continuous Control with Deep Reinforcement Learning”. Lillicrap T. P. et al.. 2016.

```
__init__(mdp_info,    policy_class,    policy_params,    actor_params,    actor_optimizer,
        critic_params,    batch_size,    initial_replay_size,    max_replay_size,    tau,    policy_delay=1,
        critic_fit_params=None)
```

Constructor.

**Parameters**

- `policy_class` ([Policy](#)) – class of the policy;
- `policy_params` (*dict*) – parameters of the policy to build;
- `actor_params` (*dict*) – parameters of the actor approximator to build;
- `actor_optimizer` (*dict*) – parameters to specify the actor optimizer algorithm;
- `critic_params` (*dict*) – parameters of the critic approximator to build;
- `batch_size` (*int*) – the number of samples in a batch;
- `initial_replay_size` (*int*) – the number of samples to collect before starting the learning;
- `max_replay_size` (*int*) – the maximum number of samples in the replay memory;
- `tau` (*float*) – value of coefficient for soft updates;
- `policy_delay` (*int*, `1`) – the number of updates of the critic after which an actor update is implemented;
- `critic_fit_params` (*dict*, `None`) – parameters of the fitting algorithm of the critic approximator;

**fit** (*dataset*)

Fit step.

**Parameters** `dataset` (*list*) – the dataset.

**\_next\_q** (*next\_state, absorbing*)

**Parameters**

- `next_state` (*np.ndarray*) – the states where next action has to be evaluated;
- `absorbing` (*np.ndarray*) – the absorbing flag for the states in `next_state`.

**Returns** Action-values returned by the critic for `next_state` and the action returned by the actor.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**\_add\_save\_attr** (*\*\*attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_optimize\_actor\_parameters** (*loss*)

Method used to update actor parameters to maximize a given loss.

**Parameters** `loss` (*torch.tensor*) – the loss computed by the algorithm.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**classmethod load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD3(mdp_info, policy_class, policy_params, actor_params, actor_optimizer, critic_params, batch_size, initial_replay_size, max_replay_size, tau, policy_delay=2, noise_std=0.2, noise_clip=0.5, critic_fit_params=None)
Bases: mushroom_rl.algorithms.actor_critic.deep_actor_critic.ddpg.DDPG
```

Twin Delayed DDPG algorithm. “Addressing Function Approximation Error in Actor-Critic Methods”. Fujimoto S. et al.. 2018.

```
__init__(mdp_info, policy_class, policy_params, actor_params, actor_optimizer, critic_params, batch_size, initial_replay_size, max_replay_size, tau, policy_delay=2, noise_std=0.2, noise_clip=0.5, critic_fit_params=None)
```

Constructor.

#### Parameters

- **policy\_class** (`Policy`) – class of the policy;
- **policy\_params** (`dict`) – parameters of the policy to build;
- **actor\_params** (`dict`) – parameters of the actor approximator to build;
- **actor\_optimizer** (`dict`) – parameters to specify the actor optimizer algorithm;
- **critic\_params** (`dict`) – parameters of the critic approximator to build;
- **batch\_size** (`int`) – the number of samples in a batch;
- **initial\_replay\_size** (`int`) – the number of samples to collect before starting the learning;
- **max\_replay\_size** (`int`) – the maximum number of samples in the replay memory;
- **tau** (`float`) – value of coefficient for soft updates;
- **policy\_delay** (`int, 2`) – the number of updates of the critic after which an actor update is implemented;
- **noise\_std** (`float, 2`) – standard deviation of the noise used for policy smoothing;
- **noise\_clip** (`float, 5`) – maximum absolute value for policy smoothing noise;
- **critic\_fit\_params** (`dict, None`) – parameters of the fitting algorithm of the critic approximator.

```
_next_q(next_state, absorbing)
```

#### Parameters

- **next\_state** (`np.ndarray`) – the states where next action has to be evaluated;
- **absorbing** (`np.ndarray`) – the absorbing flag for the states in `next_state`.

---

**Returns** Action-values returned by the critic for `next_state` and the action returned by the actor.

**\_add\_save\_attr (\*\*attr\_dict)**  
Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_optimize\_actor\_parameters (loss)**  
Method used to update actor parameters to maximize a given loss.

**Parameters** `loss` (`torch.tensor`) – the loss computed by the algorithm.

**\_post\_load()**  
This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action (state)**  
Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**  
Called by the agent when a new episode starts.

**fit (dataset)**  
Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**classmethod load (path)**  
Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save (path)**  
Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**stop()**  
Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.SAC(mdp_info, actor_mu_params, actor_sigma_params, actor_optimizer, critic_params, batch_size, initial_replay_size, max_replay_size, warmup_transitions, tau, lr_alpha, target_entropy=None, critic_fit_params=None)
```

Bases: mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic.DeepAC

Soft Actor-Critic algorithm. “Soft Actor-Critic Algorithms and Applications”. Haarnoja T. et al.. 2019.

```
__init__(mdp_info, actor_mu_params, actor_sigma_params, actor_optimizer, critic_params, batch_size, initial_replay_size, max_replay_size, warmup_transitions, tau, lr_alpha, target_entropy=None, critic_fit_params=None)
```

Constructor.

#### Parameters

- **actor\_mu\_params** (*dict*) – parameters of the actor mean approximator to build;
- **actor\_sigma\_params** (*dict*) – parameters of the actor sigm approximator to build;
- **actor\_optimizer** (*dict*) – parameters to specify the actor optimizer algorithm;
- **critic\_params** (*dict*) – parameters of the critic approximator to build;
- **batch\_size** (*int*) – the number of samples in a batch;
- **initial\_replay\_size** (*int*) – the number of samples to collect before starting the learning;
- **max\_replay\_size** (*int*) – the maximum number of samples in the replay memory;
- **warmup\_transitions** (*int*) – number of samples to accumulate in the replay memory to start the policy fitting;
- **tau** (*float*) – value of coefficient for soft updates;
- **lr\_alpha** (*float*) – Learning rate for the entropy coefficient;
- **target\_entropy** (*float, None*) – target entropy for the policy, if None a default value is computed ;
- **critic\_fit\_params** (*dict, None*) – parameters of the fitting algorithm of the critic approximator.

```
fit(dataset)
```

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

```
_add_save_attr(**attr_dict)
```

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

### `_next_q(next_state, absorbing)`

**Parameters**

- `next_state` (`np.ndarray`) – the states where next action has to be evaluated;
- `absorbing` (`np.ndarray`) – the absorbing flag for the states in `next_state`.

**Returns** Action-values returned by the critic for `next_state` and the action returned by the actor.

### `_optimize_actor_parameters(loss)`

Method used to update actor parameters to maximize a given loss.

**Parameters** `loss` (`torch.tensor`) – the loss computed by the algorithm.

### `copy()`

**Returns** A deepcopy of the agent.

### `draw_action(state)`

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

### `episode_start()`

Called by the agent when a new episode starts.

### `classmethod load(path)`

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

### `save(path)`

Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

### `stop()`

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

### `_post_load()`

This method can be overwritten to implement logic that is executed after the loading of the agent.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.TRPO(mdp_info,
                                                               policy,
                                                               critic_params,
                                                               ent_coeff=0.0,
                                                               max_kl=0.001,
                                                               lam=1.0,
                                                               n_epochs_line_search=10,
                                                               n_epochs_cg=10,
                                                               cg_damping=0.01,
                                                               cg_residual_tol=1e-10,
                                                               quiet=True,
                                                               critic_fit_params=None)
```

Bases: *mushroom\_rl.algorithms.agent.Agent*

Trust Region Policy optimization algorithm. “Trust Region Policy Optimization”. Schulman J. et al.. 2015.

```
__init__(mdp_info, policy, critic_params, ent_coeff=0.0, max_kl=0.001, lam=1.0,
        n_epochs_line_search=10, n_epochs_cg=10, cg_damping=0.01, cg_residual_tol=1e-10, quiet=True, critic_fit_params=None)
```

Constructor.

#### Parameters

- **policy** (`TorchPolicy`) – torch policy to be learned by the algorithm
- **critic\_params** (`dict`) – parameters of the critic approximator to build;
- **ent\_coeff** (`float, 0`) – coefficient for the entropy penalty;
- **max\_kl** (`float, 0.001`) – maximum kl allowed for every policy update;
- **float** (`lam`) – lambda coefficient used by generalized advantage estimation;
- **n\_epochs\_line\_search** (`int, 10`) – maximum number of iterations of the line search algorithm;
- **n\_epochs\_cg** (`int, 10`) – maximum number of iterations of the conjugate gradient algorithm;
- **cg\_damping** (`float, 1e-2`) – damping factor for the conjugate gradient algorithm;
- **cg\_residual\_tol** (`float, 1e-10`) – conjugate gradient residual tolerance;
- **quiet** (`bool, True`) – if true, the algorithm will print debug information;
- **critic\_fit\_params** (`dict, None`) – parameters of the fitting algorithm of the critic approximator.

**fit** (`dataset`)

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**\_add\_save\_attr** (`**attr_dict`)

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load** ()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()

**Returns** A deepcopy of the agent.

**draw\_action**(*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start**()

Called by the agent when a new episode starts.

**classmethod load**(*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save**(*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop**()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.actor_critic.deep_actor_critic.PPO(mdp_info,
                                                               policy, actor_optimizer,
                                                               critic_params,
                                                               n_epochs_policy,
                                                               batch_size,
                                                               eps_ppo, lam,
                                                               quiet=True,
                                                               critic_fit_params=None)
```

Bases: *mushroom\_rl.algorithms.agent.Agent*

Proximal Policy Optimization algorithm. “Proximal Policy Optimization Algorithms”. Schulman J. et al.. 2017.

```
__init__(mdp_info, policy, actor_optimizer, critic_params, n_epochs_policy, batch_size, eps_ppo,
        lam, quiet=True, critic_fit_params=None)
```

Constructor.

**Parameters**

- **policy** ([TorchPolicy](#)) – torch policy to be learned by the algorithm
- **actor\_optimizer** (*dict*) – parameters to specify the actor optimizer algorithm;
- **critic\_params** (*dict*) – parameters of the critic approximator to build;
- **n\_epochs\_policy** (*int*) – number of policy updates for every dataset;
- **batch\_size** (*int*) – size of minibatches for every optimization step
- **eps\_ppo** (*float*) – value for probability ratio clipping;
- **float** (*lam*) – lambda coefficient used by generalized advantage estimation;
- **quiet** (*bool*, *True*) – if true, the algorithm will print debug information;

- **critic\_fit\_params** (*dict, None*) – parameters of the fitting algorithm of the critic approximator.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**\_add\_save\_attr** (*\*\*attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start** ()

Called by the agent when a new episode starts.

**classmethod** **load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

## 3.3 Policy search

### 3.3.1 Policy gradient

```
class mushroom_rl.algorithms.policy_search.policy_gradient.REINFORCE(mdp_info,
                                                                     policy,
                                                                     learn-
                                                                     ing_rate,
                                                                     fea-
                                                                     tures=None)
Bases: mushroom_rl.algorithms.policy_search.policy_gradient.policy_gradient.
PolicyGradient
```

REINFORCE algorithm. “Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning”, Williams R. J.. 1992.

**\_\_init\_\_(*mdp\_info, policy, learning\_rate, features=None*)**  
Constructor.

**Parameters** **learning\_rate** (*float*) – the learning rate.

**\_compute\_gradient(*J*)**  
Return the gradient computed by the algorithm.

**Parameters** **J** (*list*) – list of the cumulative discounted rewards for each episode in the dataset.

**\_step\_update(*x, u, r*)**  
This function is called, when parsing the dataset, at each episode step.

**Parameters**

- **x** (*np.ndarray*) – the state at the current step;
- **u** (*np.ndarray*) – the action at the current step;
- **r** (*np.ndarray*) – the reward at the current step.

**\_episode\_end\_update()**

This function is called, when parsing the dataset, at the beginning of each episode. The implementation is dependent on the algorithm (e.g. REINFORCE updates some data structures).

**\_init\_update()**

This function is called, when parsing the dataset, at the beginning of each episode. The implementation is dependent on the algorithm (e.g. REINFORCE resets some data structure).

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_parse(*sample*)**

Utility to parse the sample.

**Parameters** **sample** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag. If provided, state is preprocessed with the features.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**\_update\_parameters(*J*)**

Update the parameters of the policy.

**Parameters** **J** (*list*) – list of the cumulative discounted rewards for each episode in the dataset.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action(*state*)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.





Episodic Natural Actor Critic algorithm. “A Survey on Policy Search for Robotics”, Deisenroth M. P., Neumann G., Peters J. 2013.

**`__init__(mdp_info, policy, learning_rate, features=None, critic_features=None)`**  
Constructor.

**Parameters** `critic_features` (*Features*, *None*) – features used by the critic.

**`_compute_gradient(J)`**  
Return the gradient computed by the algorithm.

**Parameters** `J` (*list*) – list of the cumulative discounted rewards for each episode in the dataset.

**`_step_update(x, u, r)`**  
This function is called, when parsing the dataset, at each episode step.

**Parameters**

- `x` (*np.ndarray*) – the state at the current step;
- `u` (*np.ndarray*) – the action at the current step;
- `r` (*np.ndarray*) – the reward at the current step.

**`_episode_end_update()`**

This function is called, when parsing the dataset, at the beginning of each episode. The implementation is dependent on the algorithm (e.g. REINFORCE updates some data structures).

**`_init_update()`**

This function is called, when parsing the dataset, at the beginning of each episode. The implementation is dependent on the algorithm (e.g. REINFORCE resets some data structure).

**`_add_save_attr(**attr_dict)`**

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**`_parse(sample)`**

Utility to parse the sample.

**Parameters** `sample` (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag. If provided, state is preprocessed with the features.

**`_post_load()`**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**`_update_parameters(J)`**

Update the parameters of the policy.

**Parameters** `J` (*list*) – list of the cumulative discounted rewards for each episode in the dataset.

**`copy()`**

**Returns** A deepcopy of the agent.

**`draw_action(state)`**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**fit(dataset)**

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

### 3.3.2 Black-Box optimization

```
class mushroom_rl.algorithms.policy_search.black_box_optimization.RWR(mdp_info,  
dis-  
tribu-  
tion,  
pol-  
icy,  
beta,  
fea-  
tures=None)
```

Bases: mushroom\_rl.algorithms.policy\_search.black\_box\_optimization.  
black\_box\_optimization.BlackBoxOptimization

Reward-Weighted Regression algorithm. “A Survey on Policy Search for Robotics”, Deisenroth M. P., Neumann G., Peters J.. 2013.

**\_\_init\_\_(mdp\_info, distribution, policy, beta, features=None)**

Constructor.

**Parameters** **beta** (*float*) – the temperature for the exponential reward transformation.

**\_update(Jep, theta)**

Function that implements the update routine of distribution parameters. Every black box algorithms should implement this function with the proper update.

**Parameters**

- **Jep** (*np.ndarray*) – a vector containing the J of the considered trajectories;
- **theta** (*np.ndarray*) – a matrix of policy parameters of the considered trajectories.

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**`_post_load()`**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**`copy()`**

**Returns** A deepcopy of the agent.

**`draw_action(state)`**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**`episode_start()`**

Called by the agent when a new episode starts.

**`fit(dataset)`**

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**`classmethod load(path)`**

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**`save(path)`**

Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**`stop()`**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**class** `mushroom_rl.algorithms.policy_search.black_box_optimization.PGPE` (`mdp_info`,  
`dis-`  
`tri-`  
`bu-`  
`tion,`  
`pol-`  
`icy,`  
`learn-`  
`ing_rate,`  
`fea-`  
`tures=None`)

Bases: `mushroom_rl.algorithms.policy_search.black_box_optimization.black_box_optimization.BlackBoxOptimization`

Policy Gradient with Parameter Exploration algorithm. “A Survey on Policy Search for Robotics”, Deisenroth M. P., Neumann G., Peters J.. 2013.

**`__init__(mdp_info, distribution, policy, learning_rate, features=None)`**  
Constructor.

**Parameters** `learning_rate` (`Parameter`) – the learning rate for the gradient step.

**`_update(Jep, theta)`**

Function that implements the update routine of distribution parameters. Every black box algorithms should implement this function with the proper update.

**Parameters**

- **Jep** (*np.ndarray*) – a vector containing the J of the considered trajectories;
- **theta** (*np.ndarray*) – a matrix of policy parameters of the considered trajectories.

**\_add\_save\_attr** (\*\**attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.policy_search.black_box_optimization.REPS (mdp_info,  
dis-  
tri-  
bu-  
tion,  
pol-  
icy,  
eps,  
fea-  
tures=None)  
Bases: mushroom_rl.algorithms.policy_search.black_box_optimization.  
black_box_optimization.BlackBoxOptimization
```

Episodic Relative Entropy Policy Search algorithm. “A Survey on Policy Search for Robotics”, Deisenroth M. P., Neumann G., Peters J.. 2013.

**\_\_init\_\_(*mdp\_info, distribution, policy, eps, features=None*)**  
Constructor.

**Parameters** **eps** (*float*) – the maximum admissible value for the Kullback-Leibler divergence between the new distribution and the previous one at each update step.

**\_update(*Jep, theta*)**

Function that implements the update routine of distribution parameters. Every black box algorithms should implement this function with the proper update.

**Parameters**

- **Jep** (*np.ndarray*) – a vector containing the J of the considered trajectories;
- **theta** (*np.ndarray*) – a matrix of policy parameters of the considered trajectories.

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action(*state*)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**fit(*dataset*)**

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load(*path*)**

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(*path*)**

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

## 3.4 Value-Based

### 3.4.1 TD

**class** mushroom\_rl.algorithms.value.td.**SARSA** (*mdp\_info*, *policy*, *learning\_rate*)

Bases: mushroom\_rl.algorithms.value.td.td.TD

SARSA algorithm.

**\_\_init\_\_** (*mdp\_info*, *policy*, *learning\_rate*)

Constructor.

#### Parameters

- **approximator** (*object*) – the approximator to use to fit the Q-function;
- **learning\_rate** (*Parameter*) – the learning rate.

**\_update** (*state*, *action*, *reward*, *next\_state*, *absorbing*)

Update the Q-table.

#### Parameters

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**\_add\_save\_attr** (\*\**attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**static \_parse** (*dataset*)

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** **dataset** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**\_post\_load**()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy**()

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start**()

Called by the agent when a new episode starts.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**class** mushroom\_rl.algorithms.value.td.**SARSA****Lambda** (*mdp\_info*, *policy*, *learning\_rate*,  
*lambda\_coeff*, *trace*=’replacing’)

Bases: mushroom\_rl.algorithms.value.td.td.TD

The SARSA(lambda) algorithm for finite MDPs.

**\_\_init\_\_** (*mdp\_info*, *policy*, *learning\_rate*, *lambda\_coeff*, *trace*=’replacing’)

Constructor.

**Parameters**

- **lambda\_coeff** (*float*) – eligibility trace coefficient;
- **trace** (*str*, ‘replacing’) – type of eligibility trace to use.

**\_update** (*state*, *action*, *reward*, *next\_state*, *absorbing*)

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**episode\_start** ()

Called by the agent when a new episode starts.

**\_add\_save\_attr** (\*\**attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**static \_parse** (*dataset*)

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** **dataset** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**`_post_load()`**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**`copy()`**

**Returns** A deepcopy of the agent.

**`draw_action(state)`**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**`fit(dataset)`**

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**`classmethod load(path)`**

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**`save(path)`**

Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**`stop()`**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**`class mushroom_rl.algorithms.value.td.ExpectedSARSA(mdp_info, policy, learning_rate)`**

Bases: `mushroom_rl.algorithms.value.td.TD`

Expected SARSA algorithm. “A theoretical and empirical analysis of Expected Sarsa”. Seijen H. V. et al.. 2009.

**`__init__(mdp_info, policy, learning_rate)`**

Constructor.

**Parameters**

- `approximator` (`object`) – the approximator to use to fit the Q-function;
- `learning_rate` (`Parameter`) – the learning rate.

**`_update(state, action, reward, next_state, absorbing)`**

Update the Q-table.

**Parameters**

- `state` (`np.ndarray`) – state;
- `action` (`np.ndarray`) – action;
- `reward` (`np.ndarray`) – reward;
- `next_state` (`np.ndarray`) – next state;
- `absorbing` (`np.ndarray`) – absorbing flag.

`_add_save_attr(**attr_dict)`

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

`static _parse(dataset)`

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** `dataset` (`list`) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

`_post_load()`

This method can be overwritten to implement logic that is executed after the loading of the agent.

`copy()`

**Returns** A deepcopy of the agent.

`draw_action(state)`

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

`episode_start()`

Called by the agent when a new episode starts.

`fit(dataset)`

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

`classmethod load(path)`

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

`save(path)`

Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

`stop()`

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

`class mushroom_rl.algorithms.value.td.QLearning(mdp_info, policy, learning_rate)`

Bases: `mushroom_rl.algorithms.value.td.td.TD`

Q-Learning algorithm. “Learning from Delayed Rewards”. Watkins C.J.C.H.. 1989.

`__init__(mdp_info, policy, learning_rate)`

Constructor.

**Parameters**

- `approximator` (`object`) – the approximator to use to fit the Q-function;
- `learning_rate` (`Parameter`) – the learning rate.

**\_update**(state, action, reward, next\_state, absorbing)

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**\_add\_save\_attr**(\*\*attr\_dict)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**static \_parse**(dataset)

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** **dataset** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**\_post\_load**()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy**()

**Returns** A deepcopy of the agent.

**draw\_action**(state)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start**()

Called by the agent when a new episode starts.

**fit**(dataset)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load**(path)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save**(path)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop**()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.value.td.DoubleQLearning(mdp_info, policy, learning_rate)
    Bases: mushroom_rl.algorithms.value.td.td.TD
    Double Q-Learning algorithm. “Double Q-Learning”. Hasselt H. V.. 2010.

__init__(mdp_info, policy, learning_rate)
    Constructor.

    Parameters
        • approximator (object) – the approximator to use to fit the Q-function;
        • learning_rate (Parameter) – the learning rate.

_update(state, action, reward, next_state, absorbing)
    Update the Q-table.

    Parameters
        • state (np.ndarray) – state;
        • action (np.ndarray) – action;
        • reward (np.ndarray) – reward;
        • next_state (np.ndarray) – next state;
        • absorbing (np.ndarray) – absorbing flag.

_add_save_attr(**attr_dict)
    Add attributes that should be saved for an agent.

    Parameters attr_dict (dict) – dictionary of attributes mapped to the method that should be used to save and load them.

static _parse(dataset)
    Utility to parse the dataset that is supposed to contain only a sample.

    Parameters dataset (list) – the current episode step.

    Returns A tuple containing state, action, reward, next state, absorbing and last flag.

_post_load()
    This method can be overwritten to implement logic that is executed after the loading of the agent.

copy()
    Returns A deepcopy of the agent.

draw_action(state)
    Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

    Parameters state (np.ndarray) – the state where the agent is.

    Returns The action to be executed.

episode_start()
    Called by the agent when a new episode starts.

fit(dataset)
    Fit step.

    Parameters dataset (list) – the dataset.

classmethod load(path)
    Load and deserialize the agent from the given location on disk.
```

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

### `save(path)`

Serialize and save the agent to the given path on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

### `stop()`

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**class** `mushroom_rl.algorithms.value.td.SpeedyQLearning` (*mdp\_info*, *policy*, *learning\_rate*)  
Bases: `mushroom_rl.algorithms.value.td.td.TD`

Speedy Q-Learning algorithm. “Speedy Q-Learning”. Ghavamzadeh et. al.. 2011.

### `__init__(mdp_info, policy, learning_rate)`

Constructor.

#### Parameters

- `approximator` (*object*) – the approximator to use to fit the Q-function;
- `learning_rate` (`Parameter`) – the learning rate.

### `_update(state, action, reward, next_state, absorbing)`

Update the Q-table.

#### Parameters

- `state` (`np.ndarray`) – state;
- `action` (`np.ndarray`) – action;
- `reward` (`np.ndarray`) – reward;
- `next_state` (`np.ndarray`) – next state;
- `absorbing` (`np.ndarray`) – absorbing flag.

### `_add_save_attr(**attr_dict)`

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

### `static _parse(dataset)`

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** `dataset` (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

### `_post_load()`

This method can be overwritten to implement logic that is executed after the loading of the agent.

### `copy()`

**Returns** A deepcopy of the agent.

### `draw_action(state)`

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**fit(dataset)**

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**class** mushroom\_rl.algorithms.value.td.RLearning(*mdp\_info*, *policy*, *learning\_rate*, *beta*)

Bases: mushroom\_rl.algorithms.value.td.TD

R-Learning algorithm. “A Reinforcement Learning Method for Maximizing Undiscounted Rewards”. Schwartz A.. 1993.

**\_\_init\_\_(mdp\_info, policy, learning\_rate, beta)**

Constructor.

**Parameters** **beta** (**Parameter**) – beta coefficient.

**\_update(state, action, reward, next\_state, absorbing)**

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**static \_parse(dataset)**

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** **dataset** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**\_post\_load()**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**

**Returns** A deepcopy of the agent.

**draw\_action(state)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start()**

Called by the agent when a new episode starts.

**fit(dataset)**

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.value.td.WeightedQLearning(mdp_info, policy, learning_rate, sampling=True, precision=1000)
```

Bases: `mushroom_rl.algorithms.value.td.TD`

Weighted Q-Learning algorithm. “Estimating the Maximum Expected Value through Gaussian Approximation”. D’Eramo C. et. al.. 2016.

**\_\_init\_\_(mdp\_info, policy, learning\_rate, sampling=True, precision=1000)**

Constructor.

**Parameters**

- **sampling** (*bool*, `True`) – use the approximated version to speed up the computation;
- **precision** (*int*, `1000`) – number of samples to use in the approximated version.

**\_update(state, action, reward, next\_state, absorbing)**

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

```
_next_q(next_state)
    Parameters next_state (np.ndarray) – the state where next action has to be evaluated.
    Returns The weighted estimator value in next_state.
_add_save_attr(**attr_dict)
    Add attributes that should be saved for an agent.
    Parameters attr_dict (dict) – dictionary of attributes mapped to the method that should
        be used to save and load them.
static _parse(dataset)
    Utility to parse the dataset that is supposed to contain only a sample.
    Parameters dataset (list) – the current episode step.
    Returns A tuple containing state, action, reward, next state, absorbing and last flag.
_post_load()
    This method can be overwritten to implement logic that is executed after the loading of the agent.
copy()
    Returns A deepcopy of the agent.
draw_action(state)
    Return the action to execute in the given state. It is the action returned by the policy or the action set by
    the algorithm (e.g. in the case of SARSA).
    Parameters state (np.ndarray) – the state where the agent is.
    Returns The action to be executed.
episode_start()
    Called by the agent when a new episode starts.
fit(dataset)
    Fit step.
    Parameters dataset (list) – the dataset.
classmethod load(path)
    Load and deserialize the agent from the given location on disk.
    Parameters path (string) – Relative or absolute path to the agents save location.
    Returns The loaded agent.
save(path)
    Serialize and save the agent to the given path on disk.
    Parameters path (string) – Relative or absolute path to the agents save location.
stop()
    Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup
    environments internals after a core learn/evaluate to enforce consistency.
class mushroom_rl.algorithms.value.td.RQLearning(mdp_info, policy, learning_rate,
    off_policy=False, beta=None,
    delta=None)
Bases: mushroom_rl.algorithms.value.td.TD
RQ-Learning algorithm. “Exploiting Structure and Uncertainty of Bellman Updates in Markov Decision Processes”. Tateo D. et al.. 2017.
```

---

**`_init_(mdp_info, policy, learning_rate, off_policy=False, beta=None, delta=None)`**  
Constructor.

#### Parameters

- **off\_policy** (`bool`, `False`) – whether to use the off policy setting or the online one;
- **beta** (`Parameter`, `None`) – beta coefficient;
- **delta** (`Parameter`, `None`) – delta coefficient.

**`_add_save_attr(**attr_dict)`**

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**`static _parse(dataset)`**

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** `dataset` (`list`) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**`_post_load()`**

This method can be overwritten to implement logic that is executed after the loading of the agent.

**`_update(state, action, reward, next_state, absorbing)`**

Update the Q-table.

#### Parameters

- **state** (`np.ndarray`) – state;
- **action** (`np.ndarray`) – action;
- **reward** (`np.ndarray`) – reward;
- **next\_state** (`np.ndarray`) – next state;
- **absorbing** (`np.ndarray`) – absorbing flag.

**`copy()`**

**Returns** A deepcopy of the agent.

**`draw_action(state)`**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**`episode_start()`**

Called by the agent when a new episode starts.

**`fit(dataset)`**

Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**`classmethod load(path)`**

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**\_next\_q** (*next\_state*)

**Parameters** **next\_state** (*np.ndarray*) – the state where next action has to be evaluated.

**Returns** The weighted estimator value in ‘next\_state’.

```
class mushroom_rl.algorithms.value.td.SARSAContinuous (mdp_info, policy, approximator, learning_rate, lambda_coeff, features, approximator_params=None)
```

Bases: mushroom\_rl.algorithms.value.td.td.TD

Continuous version of SARSA(lambda) algorithm.

**\_init\_** (*mdp\_info*, *policy*, *approximator*, *learning\_rate*, *lambda\_coeff*, *features*, *approximator\_params=None*)

Constructor.

**Parameters** **lambda\_coeff** (*float*) – eligibility trace coefficient.

**\_update** (*state*, *action*, *reward*, *next\_state*, *absorbing*)

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**episode\_start** ()

Called by the agent when a new episode starts.

**\_add\_save\_attr** (\*\**attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**static \_parse** (*dataset*)

Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** **dataset** (*list*) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**\_post\_load** ()

This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()

**Returns** A deepcopy of the agent.

**draw\_action(state)**

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**fit(dataset)**

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**classmethod load(path)**

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save(path)**

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.value.td.TrueOnlineSARSALambda(mdp_info, policy, learning_rate, lambda_coeff, features, approximator_params=None)
```

Bases: mushroom\_rl.algorithms.value.td.td.TD

True Online SARSA(lambda) with linear function approximation. “True Online TD(lambda)”. Seijen H. V. et al.. 2014.

**\_\_init\_\_(mdp\_info, policy, learning\_rate, lambda\_coeff, features, approximator\_params=None)**

Constructor.

**Parameters** **lambda\_coeff** (*float*) – eligibility trace coefficient.

**\_update(state, action, reward, next\_state, absorbing)**

Update the Q-table.

**Parameters**

- **state** (*np.ndarray*) – state;
- **action** (*np.ndarray*) – action;
- **reward** (*np.ndarray*) – reward;
- **next\_state** (*np.ndarray*) – next state;
- **absorbing** (*np.ndarray*) – absorbing flag.

**episode\_start()**

Called by the agent when a new episode starts.

**\_add\_save\_attr(\*\*attr\_dict)**

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**static** `_parse` (`dataset`)  
Utility to parse the dataset that is supposed to contain only a sample.

**Parameters** `dataset` (`list`) – the current episode step.

**Returns** A tuple containing state, action, reward, next state, absorbing and last flag.

**\_post\_load()**  
This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**  
**Returns** A deepcopy of the agent.

**draw\_action** (`state`)  
Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**fit** (`dataset`)  
Fit step.

**Parameters** `dataset` (`list`) – the dataset.

**classmethod** `load` (`path`)  
Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (`path`)  
Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**stop()**  
Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

### 3.4.2 Batch TD

```
class mushroom_rl.algorithms.value.batch_td.FQI (mdp_info, policy, approximator, n_iterations, approximator_params=None, fit_params=None, quiet=False, boosted=False)
```

Bases: `mushroom_rl.algorithms.value.batch_td.BatchTD`

Fitted Q-Iteration algorithm. “Tree-Based Batch Mode Reinforcement Learning”, Ernst D. et al.. 2005.

```
__init__ (mdp_info, policy, approximator, n_iterations, approximator_params=None, fit_params=None, quiet=False, boosted=False)
```

Constructor.

**Parameters**

- **n\_iterations** (`int`) – number of iterations to perform for training;

- **quiet** (*bool*, *False*) – whether to show the progress bar or not;
- **boosted** (*bool*, *False*) – whether to use boosted FQI or not.

**fit** (*dataset*)  
Fit loop.

**\_fit** (*x*)  
Single fit iteration.

**Parameters** **x** (*list*) – the dataset.

**\_fit\_boosted** (*x*)  
Single fit iteration for boosted FQI.

**Parameters** **x** (*list*) – the dataset.

**\_add\_save\_attr** (*\*\*attr\_dict*)  
Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load()**  
This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy** ()

**Returns** A deepcopy of the agent.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

**episode\_start** ()

Called by the agent when a new episode starts.

**classmethod** **load** (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**stop** ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

**class** mushroom\_rl.algorithms.value.batch\_td.**DoubleFQI** (*mdp\_info*, *policy*, *approximator*, *n\_iterations*, *approximator\_params=None*, *fit\_params=None*, *quiet=False*)

Bases: mushroom\_rl.algorithms.value.batch\_td.fqi.FQI

Double Fitted Q-Iteration algorithm. “Estimating the Maximum Expected Value in Continuous Reinforcement Learning Problems”. D’Eramo C. et al.. 2017.

`__init__(mdp_info, policy, approximator, n_iterations, approximator_params=None, fit_params=None, quiet=False)`  
Constructor.

**Parameters**

- **n\_iterations** (*int*) – number of iterations to perform for training;
- **quiet** (*bool*, *False*) – whether to show the progress bar or not;
- **boosted** (*bool*, *False*) – whether to use boosted FQI or not.

`_fit(x)`

Single fit iteration.

**Parameters** **x** (*list*) – the dataset.

`_add_save_attr(**attr_dict)`

Add attributes that should be saved for an agent.

**Parameters** **attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

`_fit_boosted(x)`

Single fit iteration for boosted FQI.

**Parameters** **x** (*list*) – the dataset.

`_post_load()`

This method can be overwritten to implement logic that is executed after the loading of the agent.

`copy()`

**Returns** A deepcopy of the agent.

`draw_action(state)`

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

`episode_start()`

Called by the agent when a new episode starts.

`fit(dataset)`

Fit loop.

`classmethod load(path)`

Load and deserialize the agent from the given location on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

`save(path)`

Serialize and save the agent to the given path on disk.

**Parameters** **path** (*string*) – Relative or absolute path to the agents save location.

`stop()`

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

---

```
class mushroom_rl.algorithms.value.batch_td.LSPI (mdp_info, policy, approximator_params=None, epsilon=0.01, fit_params=None, features=None)
Bases: mushroom_rl.algorithms.value.batch_td.batch_td.BatchTD

Least-Squares Policy Iteration algorithm. “Least-Squares Policy Iteration”. Lagoudakis M. G. and Parr R.. 2003.

__init__ (mdp_info, policy, approximator_params=None, epsilon=0.01, fit_params=None, features=None)
Constructor.

Parameters epsilon (float, 1e-2) – termination coefficient.

fit (dataset)
Fit step.

Parameters dataset (list) – the dataset.

_add_save_attr (**attr_dict)
Add attributes that should be saved for an agent.

Parameters attr_dict (dict) – dictionary of attributes mapped to the method that should be used to save and load them.

_post_load ()
This method can be overwritten to implement logic that is executed after the loading of the agent.

copy ()
Returns A deepcopy of the agent.

draw_action (state)
Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

Parameters state (np.ndarray) – the state where the agent is.

Returns The action to be executed.

episode_start ()
Called by the agent when a new episode starts.

classmethod load (path)
Load and deserialize the agent from the given location on disk.

Parameters path (string) – Relative or absolute path to the agents save location.

Returns The loaded agent.

save (path)
Serialize and save the agent to the given path on disk.

Parameters path (string) – Relative or absolute path to the agents save location.

stop ()
Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.
```

### 3.4.3 DQN

```
class mushroom_rl.algorithms.value.dqn.DQN(mdp_info, policy, approximator,
                                             approximator_params, batch_size,
                                             target_update_frequency, replay_memory=None, initial_replay_size=500,
                                             max_replay_size=5000, fit_params=None,
                                             n_approximators=1, clip_reward=True)
```

Bases: *mushroom\_rl.algorithms.agent.Agent*

Deep Q-Network algorithm. “Human-Level Control Through Deep Reinforcement Learning”. Mnih V. et al.. 2015.

```
__init__(mdp_info, policy, approximator, approximator_params, batch_size, target_update_frequency, replay_memory=None, initial_replay_size=500, max_replay_size=5000, fit_params=None, n_approximators=1, clip_reward=True)
```

Constructor.

#### Parameters

- **approximator** (*object*) – the approximator to use to fit the Q-function;
- **approximator\_params** (*dict*) – parameters of the approximator to build;
- **batch\_size** (*int*) – the number of samples in a batch;
- **target\_update\_frequency** (*int*) – the number of samples collected between each update of the target network;
- **replay\_memory** (*[ReplayMemory, PrioritizedReplayMemory]*, *None*) – the object of the replay memory to use; if None, a default replay memory is created;
- **initial\_replay\_size** (*int*) – the number of samples to collect before starting the learning;
- **max\_replay\_size** (*int*) – the maximum number of samples in the replay memory;
- **fit\_params** (*dict, None*) – parameters of the fitting algorithm of the approximator;
- **n\_approximators** (*int, 1*) – the number of approximator to use in AveragedDQN;
- **clip\_reward** (*bool, True*) – whether to clip the reward or not.

**fit** (*dataset*)

Fit step.

**Parameters** **dataset** (*list*) – the dataset.

**\_update\_target()**

Update the target network.

**\_next\_q** (*next\_state, absorbing*)

#### Parameters

- **next\_state** (*np.ndarray*) – the states where next action has to be evaluated;
- **absorbing** (*np.ndarray*) – the absorbing flag for the states in *next\_state*.

**Returns** Maximum action-value for each state in *next\_state*.

**draw\_action** (*state*)

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

---

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

**\_add\_save\_attr** (`**attr_dict`)  
Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (`dict`) – dictionary of attributes mapped to the method that should be used to save and load them.

**\_post\_load()**  
This method can be overwritten to implement logic that is executed after the loading of the agent.

**copy()**

**Returns** A deepcopy of the agent.

**episode\_start()**  
Called by the agent when a new episode starts.

**classmethod load** (`path`)  
Load and deserialize the agent from the given location on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

**save** (`path`)  
Serialize and save the agent to the given path on disk.

**Parameters** `path` (`string`) – Relative or absolute path to the agents save location.

**stop()**  
Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.value.dqn.DoubleDQN(mdp_info, policy, approximator, approximator_params, batch_size, target_update_frequency, replay_memory=None, initial_replay_size=500, max_replay_size=5000, fit_params=None, n_approximators=1, clip_reward=True)
```

Bases: `mushroom_rl.algorithms.value.dqn.DQN`

Double DQN algorithm. “Deep Reinforcement Learning with Double Q-Learning”. Hasselt H. V. et al.. 2016.

**\_next\_q** (`next_state, absorbing`)

**Parameters**

- `next_state` (`np.ndarray`) – the states where next action has to be evaluated;
- `absorbing` (`np.ndarray`) – the absorbing flag for the states in `next_state`.

**Returns** Maximum action-value for each state in `next_state`.

**\_\_init\_\_** (`mdp_info, policy, approximator, approximator_params, batch_size, target_update_frequency, replay_memory=None, initial_replay_size=500, max_replay_size=5000, fit_params=None, n_approximators=1, clip_reward=True`)  
Constructor.

**Parameters**

- **approximator** (*object*) – the approximator to use to fit the Q-function;
- **approximator\_params** (*dict*) – parameters of the approximator to build;
- **batch\_size** (*int*) – the number of samples in a batch;
- **target\_update\_frequency** (*int*) – the number of samples collected between each update of the target network;
- **replay\_memory** ([`ReplayMemory`, `PrioritizedReplayMemory`], `None`) – the object of the replay memory to use; if `None`, a default replay memory is created;
- **initial\_replay\_size** (*int*) – the number of samples to collect before starting the learning;
- **max\_replay\_size** (*int*) – the maximum number of samples in the replay memory;
- **fit\_params** (*dict*, `None`) – parameters of the fitting algorithm of the approximator;
- **n\_approximators** (*int*, *1*) – the number of approximator to use in AveragedDQN;
- **clip\_reward** (*bool*, `True`) – whether to clip the reward or not.

#### `_add_save_attr(**attr_dict)`

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

#### `_post_load()`

This method can be overwritten to implement logic that is executed after the loading of the agent.

#### `_update_target()`

Update the target network.

#### `copy()`

**Returns** A deepcopy of the agent.

#### `draw_action(state)`

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action to be executed.

#### `episode_start()`

Called by the agent when a new episode starts.

#### `fit(dataset)`

Fit step.

**Parameters** `dataset` (*list*) – the dataset.

#### `classmethod load(path)`

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

#### `save(path)`

Serialize and save the agent to the given path on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**stop()**

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```
class mushroom_rl.algorithms.value.dqn.AveragedDQN(mdp_info, policy, approximator,  
**params)
```

Bases: mushroom\_rl.algorithms.value.dqn.DQN

Averaged-DQN algorithm. “Averaged-DQN: Variance Reduction and Stabilization for Deep Reinforcement Learning”. Anschel O. et al.. 2017.

```
__init__(mdp_info, policy, approximator, **params)  
Constructor.
```

**Parameters**

- **approximator** (*object*) – the approximator to use to fit the Q-function;
- **approximator\_params** (*dict*) – parameters of the approximator to build;
- **batch\_size** (*int*) – the number of samples in a batch;
- **target\_update\_frequency** (*int*) – the number of samples collected between each update of the target network;
- **replay\_memory** ([*ReplayMemory*, *PrioritizedReplayMemory*], *None*) – the object of the replay memory to use; if *None*, a default replay memory is created;
- **initial\_replay\_size** (*int*) – the number of samples to collect before starting the learning;
- **max\_replay\_size** (*int*) – the maximum number of samples in the replay memory;
- **fit\_params** (*dict*, *None*) – parameters of the fitting algorithm of the approximator;
- **n\_approximators** (*int*, *1*) – the number of approximator to use in AveragedDQN;
- **clip\_reward** (*bool*, *True*) – whether to clip the reward or not.

```
_add_save_attr(**attr_dict)
```

Add attributes that should be saved for an agent.

**Parameters attr\_dict** (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

```
_post_load()
```

This method can be overwritten to implement logic that is executed after the loading of the agent.

```
copy()
```

**Returns** A deepcopy of the agent.

```
draw_action(state)
```

Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

**Parameters state** (*np.ndarray*) – the state where the agent is.

**Returns** The action to be executed.

```
episode_start()
```

Called by the agent when a new episode starts.

```
fit(dataset)
```

Fit step.

**Parameters** `dataset` (*list*) – the dataset.

**classmethod** `load` (*path*)

Load and deserialize the agent from the given location on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

**Returns** The loaded agent.

`save` (*path*)

Serialize and save the agent to the given path on disk.

**Parameters** `path` (*string*) – Relative or absolute path to the agents save location.

`stop` ()

Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

`_update_target` ()

Update the target network.

`_next_q` (*next\_state*, *absorbing*)

**Parameters**

- `next_state` (*np.ndarray*) – the states where next action has to be evaluated;
- `absorbing` (*np.ndarray*) – the absorbing flag for the states in *next\_state*.

**Returns** Maximum action-value for each state in *next\_state*.

**class** `mushroom_rl.algorithms.value.dqn.CategoricalDQN` (*mdp\_info*, *policy*, *approximator\_params*, *n\_atoms*, *v\_min*, *v\_max*, *\*\*params*)

Bases: `mushroom_rl.algorithms.value.dqn.DQN`

Categorical DQN algorithm. “A Distributional Perspective on Reinforcement Learning”. Bellemare M. et al.. 2017.

`__init__` (*mdp\_info*, *policy*, *approximator\_params*, *n\_atoms*, *v\_min*, *v\_max*, *\*\*params*)

Constructor.

**Parameters**

- `n_atoms` (*int*) – number of atoms;
- `v_min` (*float*) – minimum value of value-function;
- `v_max` (*float*) – maximum value of value-function.

`_add_save_attr` (*\*\*attr\_dict*)

Add attributes that should be saved for an agent.

**Parameters** `attr_dict` (*dict*) – dictionary of attributes mapped to the method that should be used to save and load them.

`_next_q` (*next\_state*, *absorbing*)

**Parameters**

- `next_state` (*np.ndarray*) – the states where next action has to be evaluated;
- `absorbing` (*np.ndarray*) – the absorbing flag for the states in *next\_state*.

**Returns** Maximum action-value for each state in *next\_state*.

`_post_load` ()

This method can be overwritten to implement logic that is executed after the loading of the agent.

---

```

update_target()
    Update the target network.

copy()

    Returns A deepcopy of the agent.

draw_action(state)
    Return the action to execute in the given state. It is the action returned by the policy or the action set by the algorithm (e.g. in the case of SARSA).

        Parameters state (np.ndarray) – the state where the agent is.

        Returns The action to be executed.

episode_start()
    Called by the agent when a new episode starts.

fit(dataset)
    Fit step.

        Parameters dataset (list) – the dataset.

classmethod load(path)
    Load and deserialize the agent from the given location on disk.

        Parameters path (string) – Relative or absolute path to the agents save location.

        Returns The loaded agent.

save(path)
    Serialize and save the agent to the given path on disk.

        Parameters path (string) – Relative or absolute path to the agents save location.

stop()
    Method used to stop an agent. Useful when dealing with real world environments, simulators, or to cleanup environments internals after a core learn/evaluate to enforce consistency.

```

## 3.5 Approximators

MushroomRL exposes the high-level class `Regressor` that can manage any type of function regressor. This class is a wrapper for any kind of function approximator, e.g. a scikit-learn approximator or a pytorch neural network.

### 3.5.1 Regressor

```

class mushroom_rl.approximators.regressor.Regressor(approximator,      input_shape,
                                                output_shape=(1,          ),
                                                n_actions=None,   n_models=1,
                                                **params)

```

Bases: `object`

This class implements the function to manage a function approximator. This class selects the appropriate kind of regressor to implement according to the parameters provided by the user; this makes this class the only one to use for each kind of task that has to be performed. The inference of the implementation to choose is done checking the provided values of parameters `n_actions`. If `n_actions` is provided, it means that the user wants to implement an approximator of the Q-function: if the value of `n_actions` is equal to the `output_shape` then a `QRegressor` is created, else (`output_shape` should be (1,)) an `ActionRegressor` is created.

Otherwise a GenericRegressor is created. An Ensemble model can be used for all the previous implementations listed before simply providing a `n_models` parameter greater than 1.

**`__init__(approximator, input_shape, output_shape=(1,), n_actions=None, n_models=1, **params)`**  
Constructor.

### Parameters

- **`approximator`** (*object*) – the approximator class to use to create the model;
- **`input_shape`** (*tuple*) – the shape of the input of the model;
- **`output_shape`** (*tuple, (1, )*) – the shape of the output of the model;
- **`n_actions`** (*int, None*) – number of actions considered to create a QRegressor or an ActionRegressor;
- **`n_models`** (*int, 1*) – number of models to create;
- **`**params`** (*dict*) – other parameters to create each model.

**`__call__(*, **predict_params)`**  
Call self as a function.

**`fit(*z, **fit_params)`**  
Fit the model.

### Parameters

- **`*z`** (*list*) – list of input of the model;
- **`**fit_params`** (*dict*) – parameters to use to fit the model.

**`predict(*z, **predict_params)`**  
Predict the output of the model given an input.

### Parameters

- **`*z`** (*list*) – list of input of the model;
- **`**predict_params`** (*dict*) – parameters to use to predict with the model.

**Returns** The model prediction.

### **`model`**

The model object.

**Type** Returns

### **`reset()`**

Reset the model parameters.

### **`input_shape`**

The shape of the input of the model.

**Type** Returns

### **`output_shape`**

The shape of the output of the model.

**Type** Returns

### **`weights_size`**

The shape of the weights of the model.

**Type** Returns

### **`get_weights()`**

**Returns** The weights of the model.

**set\_weights** (*w*)

**Parameters** **w** (*list*) – list of weights to be set in the model.

**diff** (\**z*)

**Parameters** **\*z** (*list*) – the input of the model.

**Returns** The derivative of the model.

### 3.5.2 Approximator

#### Linear

```
class mushroom_rl.approximators.parametric.linear.LinearApproximator(weights=None,
    in-
    put_shape=None,
    out-
    put_shape=(1,
    ),
    **kwargs)
```

Bases: object

This class implements a linear approximator.

**\_\_init\_\_** (*weights=None*, *input\_shape=None*, *output\_shape=(1, )*, *\*\*kwargs*)

Constructor.

##### Parameters

- **weights** (*np.ndarray*) – array of weights to initialize the weights of the approximator;
- **input\_shape** (*np.ndarray*, *None*) – the shape of the input of the model;
- **output\_shape** (*np.ndarray*, *(1, )*) – the shape of the output of the model;
- **\*\*kwargs** (*dict*) – other params of the approximator.

**fit** (*x*, *y*, *\*\*fit\_params*)

Fit the model.

##### Parameters

- **x** (*np.ndarray*) – input;
- **y** (*np.ndarray*) – target;
- **\*\*fit\_params** (*dict*) – other parameters used by the fit method of the regressor.

**predict** (*x*, *\*\*predict\_params*)

Predict.

##### Parameters

- **x** (*np.ndarray*) – input;
- **\*\*predict\_params** (*dict*) – other parameters used by the predict method the regressor.

**Returns** The predictions of the model.

**weights\_size**

The size of the array of weights.

**Type** Returns

**get\_weights()**

Getter.

**Returns** The set of weights of the approximator.

**set\_weights(w)**

Setter.

**Parameters** **w** (*np.ndarray*) – the set of weights to set.

**diff(state, action=None)**

Compute the derivative of the output w.r.t. state, and action if provided.

**Parameters**

- **state** (*np.ndarray*) – the state;
- **action** (*np.ndarray, None*) – the action.

**Returns** The derivative of the output w.r.t. state, and action if provided.

## Torch Approximator

```
class mushroom_rl.approximators.parametric.torch_approximator.TorchApproximator(input_shape,
                                                                                 out-
                                                                                 put_shape,
                                                                                 net-
                                                                                 work,
                                                                                 op-
                                                                                 ti-
                                                                                 mizer=None,
                                                                                 loss=None,
                                                                                 batch_size=0,
                                                                                 n_fit_targets=1
                                                                                 use_cuda=False,
                                                                                 reini-
                                                                                 tial-
                                                                                 ize=False,
                                                                                 dropout=False,
                                                                                 quiet=True,
                                                                                 **params)
```

Bases: *object*

Class to interface a pytorch model to the mushroom Regressor interface. This class implements all is needed to use a generic pytorch model and train it using a specified optimizer and objective function. This class supports also minibatches.

```
__init__(input_shape, output_shape, network, optimizer=None, loss=None, batch_size=0,
        n_fit_targets=1, use_cuda=False, reinitialize=False, dropout=False, quiet=True,
        **params)
```

Constructor.

**Parameters**

- **input\_shape** (*tuple*) – shape of the input of the network;

- **output\_shape** (*tuple*) – shape of the output of the network;
- **network** (*torch.nn.Module*) – the network class to use;
- **optimizer** (*dict*) – the optimizer used for every fit step;
- **loss** (*torch.nn.functional*) – the loss function to optimize in the fit method;
- **batch\_size** (*int*, *0*) – the size of each minibatch. If 0, the whole dataset is fed to the optimizer at each epoch;
- **n\_fit\_targets** (*int*, *1*) – the number of fit targets used by the fit method of the network;
- **use\_cuda** (*bool*, *False*) – if True, runs the network on the GPU;
- **reinitialize** (*bool*, *False*) – if True, the approximator is reinitialized at every fit call. To perform the initialization, the regressor must be initialized before;
- **method must be defined properly for the selected weights\_init** (*weights\_init*) –
- **network.** (*model*) –
- **dropout** (*bool*, *False*) – if True, dropout is applied only during train;
- **quiet** (*bool*, *True*) – if False, shows two progress bars, one for epochs and one for the minibatches;
- **params** (*dict*) – dictionary of parameters needed to construct the network.

**predict** (\*args, *output\_tensor=False*, \*\*kwargs)  
Predict.

#### Parameters

- **args** (*list*) – input;
- **output\_tensor** (*bool*, *False*) – whether to return the output as tensor or not;
- **\*\*kwargs** (*dict*) – other parameters used by the predict method the regressor.

**Returns** The predictions of the model.

**fit** (\*args, *n\_epochs=None*, *weights=None*, *epsilon=None*, *patience=1*, *validation\_split=1.0*, \*\*kwargs)  
Fit the model.

#### Parameters

- **\*args** (*list*) – input, where the last *n\_fit\_targets* elements are considered as the target, while the others are considered as input;
- **n\_epochs** (*int*, *None*) – the number of training epochs;
- **weights** (*np.ndarray*, *None*) – the weights of each sample in the computation of the loss;
- **epsilon** (*float*, *None*) – the coefficient used for early stopping;
- **patience** (*float*, *1.*) – the number of epochs to wait until stop the learning if not improving;
- **validation\_split** (*float*, *1.*) – the percentage of the dataset to use as training set;

- **\*\*kwargs** (*dict*) – other parameters used by the fit method of the regressor.

**set\_weights** (*weights*)

Setter.

**Parameters** **w** (*np.ndarray*) – the set of weights to set.

**get\_weights** ()

Getter.

**Returns** The set of weights of the approximator.

**weights\_size**

The size of the array of weights.

**Type** Returns

**diff** (\*args, \*\*kwargs)

Compute the derivative of the output w.r.t. state, and action if provided.

**Parameters**

- **state** (*np.ndarray*) – the state;
- **action** (*np.ndarray, None*) – the action.

**Returns** The derivative of the output w.r.t. state, and action if provided.

## 3.6 Distributions

**class** mushroom\_rl.distributions.distribution.**Distribution**  
Bases: object

Interface for Distributions to represent a generic probability distribution. Probability distributions are often used by black box optimization algorithms in order to perform exploration in parameter space. In literature, they are also known as high level policies.

**sample** ()

Draw a sample from the distribution.

**Returns** A random vector sampled from the distribution.

**log\_pdf** (*theta*)

Compute the logarithm of the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the log pdf is calculated

**Returns** The value of the log pdf in the specified point.

**\_\_call\_\_** (*theta*)

Compute the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the pdf is calculated

**Returns** The value of the pdf in the specified point.

**mle** (*theta, weights=None*)

Compute the (weighted) maximum likelihood estimate of the points, and update the distribution accordingly.

**Parameters**

- **theta** (*np.ndarray*) – a set of points, every row is a sample

- **weights** (*np.ndarray, None*) – a vector of weights. If specified the weighted maximum likelihood estimate is computed instead of the plain maximum likelihood. The number of elements of this vector must be equal to the number of rows of the theta matrix.

**diff\_log(theta)**

Compute the derivative of the gradient of the probability density function in the specified point.

#### Parameters

- **theta** (*np.ndarray*) – the point where the gradient of the log pdf is
- **calculated** –

**Returns** The gradient of the log pdf in the specified point.

**diff(theta)**

Compute the derivative of the probability density function, in the specified point. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_\rho p(\theta) = p(\theta) \nabla_\rho \log p(\theta)$$

#### Parameters

- **theta** (*np.ndarray*) – the point where the gradient of the pdf is
- **calculated** –

**Returns** The gradient of the pdf in the specified point.

**get\_parameters()**

Getter.

**Returns** The current distribution parameters.

**set\_parameters(rho)**

Setter.

**Parameters** **rho** (*np.ndarray*) – the vector of the new parameters to be used by the distribution

**parameters\_size**

Property.

**Returns** The size of the distribution parameters.

**\_\_init\_\_**

Initialize self. See help(type(self)) for accurate signature.

### 3.6.1 Gaussian

**class** mushroom\_rl.distributions.gaussian.**GaussianDistribution** (*mu, sigma*)

Bases: *mushroom\_rl.distributions.distribution.Distribution*

Gaussian distribution with fixed covariance matrix. The parameters vector represents only the mean.

**\_\_init\_\_** (*mu, sigma*)

Initialize self. See help(type(self)) for accurate signature.

**sample()**

Draw a sample from the distribution.

**Returns** A random vector sampled from the distribution.

**`log_pdf(theta)`**

Compute the logarithm of the probability density function in the specified point

**Parameters** `theta` (`np.ndarray`) – the point where the log pdf is calculated

**Returns** The value of the log pdf in the specified point.

**`__call__(theta)`**

Compute the probability density function in the specified point

**Parameters** `theta` (`np.ndarray`) – the point where the pdf is calculated

**Returns** The value of the pdf in the specified point.

**`mle(theta, weights=None)`**

Compute the (weighted) maximum likelihood estimate of the points, and update the distribution accordingly.

**Parameters**

- `theta` (`np.ndarray`) – a set of points, every row is a sample
- `weights` (`np.ndarray, None`) – a vector of weights. If specified the weighted maximum likelihood estimate is computed instead of the plain maximum likelihood. The number of elements of this vector must be equal to the number of rows of the theta matrix.

**`diff_log(theta)`**

Compute the derivative of the gradient of the probability density function in the specified point.

**Parameters**

- `theta` (`np.ndarray`) – the point where the gradient of the log pdf is
- `calculated` –

**Returns** The gradient of the log pdf in the specified point.

**`get_parameters()`**

Getter.

**Returns** The current distribution parameters.

**`set_parameters(rho)`**

Setter.

**Parameters** `rho` (`np.ndarray`) – the vector of the new parameters to be used by the distribution

**`parameters_size`**

Property.

**Returns** The size of the distribution parameters.

**`diff(theta)`**

Compute the derivative of the probability density function, in the specified point. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_\rho p(\theta) = p(\theta) \nabla_\rho \log p(\theta)$$

**Parameters**

- `theta` (`np.ndarray`) – the point where the gradient of the pdf is
- `calculated` –

**Returns** The gradient of the pdf in the specified point.

```
class mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution(mu,  
std)
```

Bases: *mushroom\_rl.distributions.distribution.Distribution*

Gaussian distribution with diagonal covariance matrix. The parameters vector represents the mean and the standard deviation for each dimension.

```
__init__(mu, std)
```

Initialize self. See help(type(self)) for accurate signature.

```
sample()
```

Draw a sample from the distribution.

**Returns** A random vector sampled from the distribution.

```
log_pdf(theta)
```

Compute the logarithm of the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the log pdf is calculated

**Returns** The value of the log pdf in the specified point.

```
__call__(theta)
```

Compute the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the pdf is calculated

**Returns** The value of the pdf in the specified point.

```
mle(theta, weights=None)
```

Compute the (weighted) maximum likelihood estimate of the points, and update the distribution accordingly.

#### Parameters

- **theta** (*np.ndarray*) – a set of points, every row is a sample
- **weights** (*np.ndarray, None*) – a vector of weights. If specified the weighted maximum likelihood estimate is computed instead of the plain maximum likelihood. The number of elements of this vector must be equal to the number of rows of the theta matrix.

```
diff_log(theta)
```

Compute the derivative of the gradient of the probability density function in the specified point.

#### Parameters

- **theta** (*np.ndarray*) – the point where the gradient of the log pdf is
- **calculated** –

**Returns** The gradient of the log pdf in the specified point.

```
get_parameters()
```

Getter.

**Returns** The current distribution parameters.

```
set_parameters(rho)
```

Setter.

**Parameters** **rho** (*np.ndarray*) – the vector of the new parameters to be used by the distribution

```
parameters_size
```

Property.

**Returns** The size of the distribution parameters.

**diff(theta)**

Compute the derivative of the probability density function, in the specified point. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_\theta p(\theta) = p(\theta) \nabla_\theta \log p(\theta)$$

**Parameters**

- **theta** (*np.ndarray*) – the point where the gradient of the pdf is
- **calculated.** –

**Returns** The gradient of the pdf in the specified point.

```
class mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution(mu,
                                                                     sigma)
```

Bases: *mushroom\_rl.distributions.distribution.Distribution*

Gaussian distribution with full covariance matrix. The parameters vector represents the mean and the Cholesky decomposition of the covariance matrix. This parametrization enforce the covariance matrix to be positive definite.

**\_\_init\_\_(mu, sigma)**

Initialize self. See help(type(self)) for accurate signature.

**sample()**

Draw a sample from the distribution.

**Returns** A random vector sampled from the distribution.

**log\_pdf(theta)**

Compute the logarithm of the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the log pdf is calculated

**Returns** The value of the log pdf in the specified point.

**\_\_call\_\_(theta)**

Compute the probability density function in the specified point

**Parameters** **theta** (*np.ndarray*) – the point where the pdf is calculated

**Returns** The value of the pdf in the specified point.

**mle(theta, weights=None)**

Compute the (weighted) maximum likelihood estimate of the points, and update the distribution accordingly.

**Parameters**

- **theta** (*np.ndarray*) – a set of points, every row is a sample
- **weights** (*np.ndarray, None*) – a vector of weights. If specified the weighted maximum likelihood estimate is computed instead of the plain maximum likelihood. The number of elements of this vector must be equal to the number of rows of the theta matrix.

**diff\_log(theta)**

Compute the derivative of the gradient of the probability density function in the specified point.

**Parameters**

- **theta** (*np.ndarray*) – the point where the gradient of the log pdf is

- **calculated** –

**Returns** The gradient of the log pdf in the specified point.

### **get\_parameters()**

Getter.

**Returns** The current distribution parameters.

### **set\_parameters(rho)**

Setter.

**Parameters** **rho** (*np.ndarray*) – the vector of the new parameters to be used by the distribution

### **parameters\_size**

Property.

**Returns** The size of the distribution parameters.

### **diff(theta)**

Compute the derivative of the probability density function, in the specified point. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_\theta p(\theta) = p(\theta) \nabla_\theta \log p(\theta)$$

### **Parameters**

- **theta** (*np.ndarray*) – the point where the gradient of the pdf is
- **calculated**. –

**Returns** The gradient of the pdf in the specified point.

## 3.7 Environments

In mushroom\_rl we distinguish between two different types of environment classes:

- proper environments
- generators

While environments directly implement the `Environment` interface, generators are a set of methods used to generate finite markov chains that represent a specific environment e.g., grid worlds.

### 3.7.1 Environments

#### Atari

##### **class** `mushroom_rl.environments.atari.MaxAndSkip(env, skip, max_pooling=True)`

Bases: `gym.core.Wrapper`

##### **\_\_init\_\_(env, skip, max\_pooling=True)**

Initialize self. See help(type(self)) for accurate signature.

##### **step(action)**

Run one timestep of the environment's dynamics. When end of episode is reached, you are responsible for calling `reset()` to reset this environment's state.

Accepts an action and returns a tuple (observation, reward, done, info).

**Parameters** `action (object)` – an action provided by the agent

**Returns** agent’s observation of the current environment reward (float) : amount of reward returned after previous action done (bool): whether the episode has ended, in which case further step() calls will return undefined results info (dict): contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

**Return type** observation (object)

**reset (\*\*kwargs)**

Resets the state of the environment and returns an initial observation.

**Returns** the initial observation.

**Return type** observation (object)

**close ()**

Override close in your subclass to perform any necessary cleanup.

Environments will automatically close() themselves when garbage collected or when the program exits.

**render (mode='human', \*\*kwargs)**

Renders the environment.

The set of supported modes varies per environment. (And some environments do not support rendering at all.) By convention, if mode is:

- human: render to the current display or terminal and return nothing. Usually for human consumption.
- rgb\_array: Return an numpy.ndarray with shape (x, y, 3), representing RGB values for an x-by-y pixel image, suitable for turning into a video.
- ansi: Return a string (str) or StringIO.StringIO containing a terminal-style text representation. The text can include newlines and ANSI escape sequences (e.g. for colors).

---

**Note:**

**Make sure that your class’s metadata ‘render.modes’ key includes** the list of supported modes. It’s recommended to call super() in implementations to use the functionality of this method.

---

**Parameters** `mode (str)` – the mode to render with

Example:

```
class MyEnv(Env): metadata = {‘render.modes’: [‘human’, ‘rgb_array’]}

def render(self, mode=’human’):

    if mode == ‘rgb_array’: return np.array(...) # return RGB frame suitable for video
    elif mode == ‘human’: ... # pop up a window and render
    else: super(MyEnv, self).render(mode=mode) # just raise an exception
```

**seed (seed=None)**

Sets the seed for this env’s random number generator(s).

---

**Note:** Some environments use multiple pseudorandom number generators. We want to capture all such seeds used in order to ensure that there aren’t accidental correlations between multiple generators.

---

**Returns**

**Returns the list of seeds used in this env's random** number generators. The first value in the list should be the “main” seed, or the value which a reproducer should pass to ‘seed’. Often, the main seed equals the provided ‘seed’, but this won't be true if seed=None, for example.

**Return type** list<bigint>

**unwrapped**

Completely unwrap this env.

**Returns** The base non-wrapped gym.Env instance

**Return type** gym.Env

```
class mushroom_rl.environments.atari.LazyFrames(frames, history_length)
```

Bases: object

From OpenAI Baseline. [https://github.com/openai/baselines/blob/master/baselines/common/atari\\_wrappers.py](https://github.com/openai/baselines/blob/master/baselines/common/atari_wrappers.py)

```
__init__(frames, history_length)
```

Initialize self. See help(type(self)) for accurate signature.

```
class mushroom_rl.environments.atari.Atari(name, width=84, height=84,
                                             ends_at_life=False, max_pooling=True,
                                             history_length=4, max_no_op_actions=30)
```

Bases: *mushroom\_rl.environments.environment.Environment*

The Atari environment as presented in: “Human-level control through deep reinforcement learning”. Mnih et. al.. 2015.

```
__init__(name, width=84, height=84, ends_at_life=False, max_pooling=True, history_length=4,
        max_no_op_actions=30)
```

Constructor.

**Parameters**

- **name** (str) – id name of the Atari game in Gym;
- **width** (int, 84) – width of the screen;
- **height** (int, 84) – height of the screen;
- **ends\_at\_life** (bool, False) – whether the episode ends when a life is lost or not;
- **max\_pooling** (bool, True) – whether to do max-pooling or average-pooling of the last two frames when using NoFrameskip;
- **history\_length** (int, 4) – number of frames to form a state;
- **max\_no\_op\_actions** (int, 30) – maximum number of no-op action to execute at the beginning of an episode.

```
reset(state=None)
```

Reset the current state.

**Parameters** **state** (np.ndarray, None) – the state to set to the current state.

**Returns** The current state.

```
step(action)
```

Move the agent from its current state according to the action.

**Parameters** **action** (np.ndarray) – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**stop()**

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound(x, min\_value, max\_value)**

Method used to bound state and action variables.

**Parameters**

- `x` – the variable to bound;
- `min_value` – the minimum value;
- `max_value` – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**seed(seed)**

Set the seed of the environment.

**Parameters** `seed(float)` – the value of the seed.

**set\_episode\_end(ends\_at\_life)**

Setter.

**Parameters** `ends_at_life(bool)` – whether the episode ends when a life is lost or not.

## Car on hill

**class** `mushroom_rl.environments.car_on_hill.CarOnHill(horizon=100, gamma=0.95)`

Bases: `mushroom_rl.environments.environment.Environment`

The Car On Hill environment as presented in: “Tree-Based Batch Mode Reinforcement Learning”. Ernst D. et al.. 2005.

**\_\_init\_\_(horizon=100, gamma=0.95)**

Constructor.

**reset(state=None)**

Reset the current state.

**Parameters** `state(np.ndarray, None)` – the state to set to the current state.

**Returns** The current state.

**step(action)**

Move the agent from its current state according to the action.

**Parameters** `action(np.ndarray)` – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**static \_bound**(*x, min\_value, max\_value*)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns**seed**(*seed*)

Set the seed of the environment.

**Parameters** **seed**(*float*) – the value of the seed.

**stop**()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## DeepMind Control Suite

```
class mushroom_rl.environments.dm_control_env.DMControl(domain_name, task_name,
horizon, gamma,
task_kwargs=None,
dt=0.01,
width_screen=480,
height_screen=480,
camera_id=0)
```

Bases: *mushroom\_rl.environments.environment.Environment*

Interface for dm\_control suite Mujoco environments. It makes it possible to use every dm\_control suite Mujoco environment just providing the necessary information.

```
__init__(domain_name, task_name, horizon, gamma, task_kwargs=None, dt=0.01,
width_screen=480, height_screen=480, camera_id=0)
```

Constructor.

**Parameters**

- **domain\_name** (*str*) – name of the environment;
- **task\_name** (*str*) – name of the task of the environment;
- **horizon** (*int*) – the horizon;
- **gamma** (*float*) – the discount factor;
- **task\_kwargs** (*dict, None*) – parameters of the task;
- **dt** (*float, 0.01*) – duration of a control step;
- **width\_screen** (*int, 480*) – width of the screen;
- **height\_screen** (*int, 480*) – height of the screen;
- **camera\_id** (*int, 0*) – position of camera to render the environment;

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray, None*) – the state to set to the current state.

**Returns** The current state.

**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound** (*x, min\_value, max\_value*)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

## Finite MDP

**class** mushroom\_rl.environments.finite\_mdp.**FiniteMDP** (*p, rew, mu=None, gamma=0.9, horizon=inf*)

Bases: *mushroom\_rl.environments.environment.Environment*

Finite Markov Decision Process.

**\_\_init\_\_** (*p, rew, mu=None, gamma=0.9, horizon=inf*)

Constructor.

**Parameters**

- **p** (*np.ndarray*) – transition probability matrix;
- **rew** (*np.ndarray*) – reward matrix;
- **mu** (*np.ndarray, None*) – initial state probability distribution;
- **gamma** (*float, 0.9*) – discount factor;
- **horizon** (*int, np.inf*) – the horizon.

**reset**(state=None)

Reset the current state.

**Parameters** **state**(*np.ndarray*, *None*) – the state to set to the current state.

**Returns** The current state.

**step**(action)

Move the agent from its current state according to the action.

**Parameters** **action**(*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing **action** in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**static \_bound**(x, min\_value, max\_value)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**seed**(seed)

Set the seed of the environment.

**Parameters** **seed**(*float*) – the value of the seed.

**stop**()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## Grid World

```
class mushroom_rl.environments.grid_world.AbstractGridWorld(mdp_info, height, width, start, goal)
```

Bases: *mushroom\_rl.environments.environment.Environment*

Abstract class to build a grid world.

**\_\_init\_\_**(mdp\_info, height, width, start, goal)

Constructor.

**Parameters**

- **height**(*int*) – height of the grid;
- **width**(*int*) – width of the grid;
- **start**(*tuple*) – x-y coordinates of the goal;
- **goal**(*tuple*) – x-y coordinates of the goal.

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray, None*) – the state to set to the current state.

**Returns** The current state.

**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**static \_bound** (*x, min\_value, max\_value*)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**class** mushroom\_rl.environments.grid\_world.**GridWorld** (*height, width, goal, start=(0, 0)*)

Bases: *mushroom\_rl.environments.grid\_world.AbstractGridWorld*

Standard grid world.

**\_\_init\_\_** (*height, width, goal, start=(0, 0)*)

Constructor.

**Parameters**

- **height** (*int*) – height of the grid;
- **width** (*int*) – width of the grid;
- **start** (*tuple*) – x-y coordinates of the goal;
- **goal** (*tuple*) – x-y coordinates of the goal.

**static \_bound** (*x, min\_value, max\_value*)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;

- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

#### **info**

An object containing the info of the environment.

**Type** Returns

#### **reset (state=None)**

Reset the current state.

**Parameters** **state** (*np.ndarray*, *None*) – the state to set to the current state.

**Returns** The current state.

#### **seed (seed)**

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

#### **step (action)**

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

#### **stop ()**

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

```
class mushroom_rl.environments.grid_world.GridWorldVanHasselt (height=3,  
                                width=3,  
                                goal=(0, 2),  
                                start=(2, 0))
```

Bases: *mushroom\_rl.environments.grid\_world.AbstractGridWorld*

A variant of the grid world as presented in: “Double Q-Learning”. Hasselt H. V.. 2010.

#### **\_\_init\_\_ (height=3, width=3, goal=(0, 2), start=(2, 0))**

Constructor.

##### **Parameters**

- **height** (*int*) – height of the grid;
- **width** (*int*) – width of the grid;
- **start** (*tuple*) – x-y coordinates of the goal;
- **goal** (*tuple*) – x-y coordinates of the goal.

#### **static \_bound (x, min\_value, max\_value)**

Method used to bound state and action variables.

##### **Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray, None*) – the state to set to the current state.

**Returns** The current state.

**seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## Gym

**class** *mushroom\_rl.environments.gym\_env.Gym* (*name, horizon, gamma*)  
Bases: *mushroom\_rl.environments.environment.Environment*

Interface for OpenAI Gym environments. It makes it possible to use every Gym environment just providing the id, except for the Atari games that are managed in a separate class.

**\_\_init\_\_** (*name, horizon, gamma*)

Constructor.

**Parameters**

- **name** (*str*) – gym id of the environment;
- **horizon** (*int*) – the horizon;
- **gamma** (*float*) – the discount factor.

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray, None*) – the state to set to the current state.

**Returns** The current state.

**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

### `static _bound(x, min_value, max_value)`

Method used to bound state and action variables.

#### Parameters

- `x` – the variable to bound;
- `min_value` – the minimum value;
- `max_value` – the maximum value;

**Returns** The bounded variable.

### `info`

An object containing the info of the environment.

#### Type

Returns

### `seed(seed)`

Set the seed of the environment.

**Parameters** `seed(float)` – the value of the seed.

### `stop()`

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## Inverted pendulum

```
class mushroom_rl.environments.inverted_pendulum.InvertedPendulum(random_start=False,
                                                                    m=1.0,
                                                                    l=1.0,
                                                                    g=9.8,
                                                                    mu=0.01,
                                                                    max_u=5.0,
                                                                    horizon=5000,
                                                                    zon=5000,
                                                                    gamma=0.99)
```

Bases: `mushroom_rl.environments.environment.Environment`

The Inverted Pendulum environment (continuous version) as presented in: “Reinforcement Learning In Continuous Time and Space”. Doya K.. 2000. “Off-Policy Actor-Critic”. Degris T. et al.. 2012. “Deterministic Policy Gradient Algorithms”. Silver D. et al. 2014.

```
__init__(random_start=False, m=1.0, l=1.0, g=9.8, mu=0.01, max_u=5.0, horizon=5000,
        gamma=0.99)
```

Constructor.

#### Parameters

- `random_start (bool, False)` – whether to start from a random position or from the horizontal one;
- `m (float, 1.0)` – mass of the pendulum;
- `l (float, 1.0)` – length of the pendulum;
- `g (float, 9.8)` – gravity acceleration constant;

- **mu** (*float*,  $1e-2$ ) – friction constant of the pendulum;
- **max\_u** (*float*,  $5.0$ ) – maximum allowed input torque;
- **horizon** (*int*,  $5000$ ) – horizon of the problem;
- **gamma** (*int*,  $99$ ) – discount factor.

**reset** (*state=None*)  
Reset the current state.

**Parameters** **state** (*np.ndarray*, *None*) – the state to set to the current state.

**Returns** The current state.

**step** (*action*)  
Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**stop** ()  
Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound** (*x, min\_value, max\_value*)  
Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**  
An object containing the info of the environment.

**Type** Returns

**seed** (*seed*)  
Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

## Cart Pole

```
class mushroom_rl.environments.cart_pole.CartPole (m=2.0, M=8.0, l=0.5, g=9.8,
mu=0.01, max_u=50.0,
noise_u=10.0, horizon=3000,
gamma=0.95)
```

Bases: *mushroom\_rl.environments.environment.Environment*

The Inverted Pendulum on a Cart environment as presented in: “Least-Squares Policy Iteration”. Lagoudakis M. G. and Parr R.. 2003.

```
__init__ (m=2.0, M=8.0, l=0.5, g=9.8, mu=0.01, max_u=50.0, noise_u=10.0, horizon=3000,
gamma=0.95)
```

Constructor.

**Parameters**

- **m** (*float*, *2.0*) – mass of the pendulum;
- **M** (*float*, *8.0*) – mass of the cart;
- **l** (*float*, *5*) – length of the pendulum;
- **g** (*float*, *9.8*) – gravity acceleration constant;
- **mu** (*float*, *1e-2*) – friction constant of the pendulum;
- **max\_u** (*float*, *50.*) – maximum allowed input torque;
- **noise\_u** (*float*, *10.*) – maximum noise on the action;
- **horizon** (*int*, *3000*) – horizon of the problem;
- **gamma** (*int*, *95*) – discount factor.

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray*, *None*) – the state to set to the current state.**Returns** The current state.**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound** (*x, min\_value, max\_value*)

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.**info**

An object containing the info of the environment.

**Type** Returns**seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

## LQR

```
class mushroom_rl.environments.lqr.LQR(A, B, Q, R, max_pos=inf, max_action=inf, random_init=False, episodic=False, gamma=0.9, horizon=50)
```

Bases: *mushroom\_rl.environments.environment.Environment*

This class implements a Linear-Quadratic Regulator. This task aims to minimize the undesired deviations from nominal values of some controller settings in control problems. The system equations in this task are:

$$x_{t+1} = Ax_t + Bu_t$$

where *x* is the state and *u* is the control signal.

The reward function is given by:

$$r_t = - (x_t^T Q x_t + u_t^T R u_t)$$

“Policy gradient approaches for multi-objective sequential decision making”. Parisi S., Pirotta M., Smacchia N., Bascetta L., Restelli M.. 2014

```
__init__(A, B, Q, R, max_pos=inf, max_action=inf, random_init=False, episodic=False, gamma=0.9, horizon=50)
```

Constructor.

**Args:** *A* (np.ndarray): the state dynamics matrix; *B* (np.ndarray): the action dynamics matrix; *Q* (np.ndarray): reward weight matrix for state; *R* (np.ndarray): reward weight matrix for action; *max\_pos* (float, np.inf): maximum value of the state; *max\_action* (float, np.inf): maximum value of the action; *random\_init* (bool, False): start from a random state; *episodic* (bool, False): end the episode when the state goes over the threshold; *gamma* (float, 0.9): discount factor; *horizon* (int, 50): horizon of the mdp.

```
static generate(dimensions, max_pos=inf, max_action=inf, eps=0.1, index=0, random_init=False, episodic=False, gamma=0.9, horizon=50)
```

Factory method that generates an lqr with identity dynamics and symmetric reward matrices.

### Parameters

- **dimensions** (int) – number of state-action dimensions;
- **max\_pos** (float, np.inf) – maximum value of the state;
- **max\_action** (float, np.inf) – maximum value of the action;
- **eps** (double, 1) – reward matrix weights specifier;
- **index** (int, 0) – selector for the principal state;
- **random\_init** (bool, False) – start from a random state;
- **episodic** (bool, False) – end the episode when the state goes over the threshold;
- **gamma** (float, 0.9) – discount factor;
- **horizon** (int, 50) – horizon of the mdp.

```
reset(state=None)
```

Reset the current state.

**Parameters** **state** (np.ndarray, None) – the state to set to the current state.

**Returns** The current state.

```
step(action)
```

Move the agent from its current state according to the action.

**Parameters** `action` (`np.ndarray`) – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**static \_bound** (`x, min_value, max_value`)

Method used to bound state and action variables.

#### Parameters

- `x` – the variable to bound;
- `min_value` – the minimum value;
- `max_value` – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**seed** (`seed`)

Set the seed of the environment.

**Parameters** `seed` (`float`) – the value of the seed.

**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## Mujoco

**class** `mushroom_rl.environments.mujoco.ObservationType`

Bases: `enum.Enum`

An enum indicating the type of data that should be added to the observation of the environment, can be Joint-/Body-/Site- positions and velocities.

**class** `mushroom_rl.environments.mujoco.MuJoCo` (`file_name, actuation_spec, observation_spec, gamma, horizon, n_substeps=1, n_intermediate_steps=1, additional_data_spec=None, collision_groups=None`)

Bases: `mushroom_rl.environments.environment.Environment`

Class to create a Mushroom environment using the MuJoCo simulator.

**\_\_init\_\_** (`file_name, actuation_spec, observation_spec, gamma, horizon, n_substeps=1, n_intermediate_steps=1, additional_data_spec=None, collision_groups=None`)

Constructor.

#### Parameters

- `file_name` (`string`) – The path to the XML file with which the environment should be created;
- `actuation_spec` (`list`) – A list specifying the names of the joints which should be controllable by the agent. Can be left empty when all actuators should be used;

- **observation\_spec** (*list*) – A list containing the names of data that should be made available to the agent as an observation and their type (ObservationType). An entry in the list is given by: (name, type);
- **gamma** (*float*) – The discounting factor of the environment;
- **horizon** (*int*) – The maximum horizon for the environment;
- **n\_substeps** (*int*) – The number of substeps to use by the MuJoCo simulator. An action given by the agent will be applied for n\_substeps before the agent receives the next observation and can act accordingly;
- **n\_intermediate\_steps** (*int*) – The number of steps between every action taken by the agent. Similar to n\_substeps but allows the user to modify, control and access intermediate states.
- **additional\_data\_spec** (*list*) – A list containing the data fields of interest, which should be read from or written to during simulation. The entries are given as the following tuples: (key, name, type) key is a string for later referencing in the “read\_data” and “write\_data” methods. The name is the name of the object in the XML specification and the type is the ObservationType;
- **collision\_groups** (*list*) – A list containing groups of geoms for which collisions should be checked during simulation via check\_collision. The entries are given as: (key, geom\_names), where key is a string for later referencing in the “check\_collision” method, and geom\_names is a list of geom names in the XML specification.

**seed** (*seed*)

Set the seed of the environment.

**Parameters** **seed** (*float*) – the value of the seed.

**reset** (*state=None*)

Reset the current state.

**Parameters** **state** (*np.ndarray, None*) – the state to set to the current state.

**Returns** The current state.

**stop** ()

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**step** (*action*)

Move the agent from its current state according to the action.

**Parameters** **action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**\_preprocess\_action** (*action*)

Compute a transformation of the action provided to the environment.

**Parameters** **action** (*np.ndarray*) – numpy array with the actions provided to the environment.

**Returns** The action to be used for the current step

**\_step\_init** (*state, action*)

Allows information to be initialized at the start of a step.

**\_compute\_action (action)**

Compute a transformation of the action at every intermediate step. Useful to add control signals simulated directly in python.

**Parameters** **action** (*np.ndarray*) – numpy array with the actions provided at every step.

**Returns** The action to be set in the actual mujoco simulation.

**\_simulation\_pre\_step ()**

**Allows information to be accesed and changed at every intermediate step** before taking a step in the mujoco simulation. Can be usefull to apply an external force/torque to the specified bodies.

**ex: apply a force over X to the torso:** force = [200, 0, 0] torque = [0, 0, 0]  
self.sim.data.xfrc\_applied[self.sim.model.\_body\_name2id[“torso”],:] = force + torque

**\_simulation\_post\_step ()**

**Allows information to be accesed at every intermediate step** after taking a step in the mujoco simulation. Can be usefull to average forces over all intermediate steps.

**\_step\_finalize ()**

Allows information to be accesed at the end of a step.

**read\_data (name)**

Read data form the MuJoCo data structure.

**Parameters** **name** (*string*) – A name referring to an entry contained the additional\_data\_spec list handed to the constructor.

**Returns** The desired data as a one-dimensional numpy array.

**write\_data (name, value)**

Write data to the MuJoCo data structure.

**Parameters**

- **name** (*string*) – A name referring to an entry contained in the additional\_data\_spec list handed to the constructor;
- **value** (*ndarray*) – The data that should be written.

**check\_collision (group1, group2)**

Check for collision between the specified groups.

**Parameters**

- **group1** (*string*) – A name referring to an entry contained in the collision\_groups list handed to the constructor;
- **group2** (*string*) – A name referring to an entry contained in the collision\_groups list handed to the constructor.

**Returns** A flag indicating whether a collision occurred between the given groups or not.

**get\_collision\_force (group1, group2)**

Returns the collision force and torques between the specified groups.

**Parameters**

- **group1** (*string*) – A name referring to an entry contained in the collision\_groups list handed to the constructor;
- **group2** (*string*) – A name referring to an entry contained in the collision\_groups list handed to the constructor.

**Returns** A 6D vector specifying the collision forces/torques[3D force + 3D torque] between the given groups. Vector of 0's in case there was no collision. <http://mujoco.org/book/programming.html#siContact>

**reward**(*state, action, next\_state*)

Compute the reward based on the given transition.

#### Parameters

- **state** (*np.array*) – the current state of the system;
- **action** (*np.array*) – the action that is applied in the current state;
- **next\_state** (*np.array*) – the state reached after applying the given action.

**Returns** The reward as a floating point scalar value.

**is\_absorbing**(*state*)

Check whether the given state is an absorbing state or not.

**Parameters** **state** (*np.array*) – the state of the system.

**Returns** A boolean flag indicating whether this state is absorbing or not.

**static \_bound**(*x, min\_value, max\_value*)

Method used to bound state and action variables.

#### Parameters

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**

An object containing the info of the environment.

**Type** Returns

**setup()**

A function that allows to execute setup code after an environment reset.

## Puddle World

```
class mushroom_rl.environments.puddle_world.PuddleWorld(start=None, goal=None,
goal_threshold=0.1,
noise_step=0.025,
noise_reward=0,
reward_goal=0.0,
thrust=0.05, puddle_center=None,
puddle_width=None,
gamma=0.99, horizon=5000)
```

Bases: *mushroom\_rl.environments.environment.Environment*

Puddle world as presented in: “Off-Policy Actor-Critic”. Degris T. et al.. 2012.

---

**\_\_init\_\_(start=None, goal=None, goal\_threshold=0.1, noise\_step=0.025, noise\_reward=0, reward\_goal=0.0, thrust=0.05, puddle\_center=None, puddle\_width=None, gamma=0.99, horizon=5000)**  
Constructor.

**Parameters**

- **start** (*np.array*, *None*) – starting position of the agent;
- **goal** (*np.array*, *None*) – goal position;
- **goal\_threshold** (*float*, *1*) – distance threshold of the agent from the goal to consider it reached;
- **noise\_step** (*float*, *0.025*) – noise in actions;
- **noise\_reward** (*float*, *0*) – standard deviation of gaussian noise in reward;
- **reward\_goal** (*float*, *0*) – reward obtained reaching goal state;
- **thrust** (*float*, *0.05*) – distance walked during each action;
- **puddle\_center** (*np.array*, *None*) – center of the puddle;
- **puddle\_width** (*np.array*, *None*) – width of the puddle;

**reset(state=None)**

Reset the current state.

**Parameters state** (*np.ndarray*, *None*) – the state to set to the current state.**Returns** The current state.**step(action)**

Move the agent from its current state according to the action.

**Parameters action** (*np.ndarray*) – the action to execute.**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).**stop()**

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound(x, min\_value, max\_value)**

Method used to bound state and action variables.

**Parameters**

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.**info**

An object containing the info of the environment.

**Type** Returns**seed(seed)**

Set the seed of the environment.

**Parameters seed** (*float*) – the value of the seed.

## Segway

```
class mushroom_rl.environments.segway.Segway(random_start=False)
Bases: mushroom_rl.environments.environment.Environment
```

The Segway environment (continuous version) as presented in: “Deep Learning for Actor-Critic Reinforcement Learning”. Xueli Jia. 2015.

```
__init__(random_start=False)
```

Constructor.

**Parameters** `random_start` (`bool`, `False`) – whether to start from a random position or from the horizontal one.

```
reset(state=None)
```

Reset the current state.

**Parameters** `state` (`np.ndarray`, `None`) – the state to set to the current state.

**Returns** The current state.

```
step(action)
```

Move the agent from its current state according to the action.

**Parameters** `action` (`np.ndarray`) – the action to execute.

**Returns** The state reached by the agent executing `action` in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

```
static _bound(x, min_value, max_value)
```

Method used to bound state and action variables.

**Parameters**

- `x` – the variable to bound;
- `min_value` – the minimum value;
- `max_value` – the maximum value;

**Returns** The bounded variable.

```
info
```

An object containing the info of the environment.

**Type** Returns

```
seed(seed)
```

Set the seed of the environment.

**Parameters** `seed` (`float`) – the value of the seed.

```
stop()
```

Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

## Ship steering

```
class mushroom_rl.environments.ship_steering.ShipSteering(small=True,
                                                               n_steps_action=3)
Bases: mushroom_rl.environments.environment.Environment
```

The Ship Steering environment as presented in: “Hierarchical Policy Gradient Algorithms”. Ghavamzadeh M. and Mahadevan S.. 2013.

**\_\_init\_\_(small=True, n\_steps\_action=3)**  
Constructor.

#### Parameters

- **small** (*bool*, *True*) – whether to use a small state space or not.
- **n\_steps\_action** (*int*, *3*) – number of integration intervals for each step of the mdp.

**reset(state=None)**  
Reset the current state.

**Parameters state** (*np.ndarray*, *None*) – the state to set to the current state.

**Returns** The current state.

**step(action)**  
Move the agent from its current state according to the action.

**Parameters action** (*np.ndarray*) – the action to execute.

**Returns** The state reached by the agent executing *action* in its current state, the reward obtained in the transition and a flag to signal if the next state is absorbing. Also an additional dictionary is returned (possibly empty).

**stop()**  
Method used to stop an mdp. Useful when dealing with real world environments, simulators, or when using openai-gym rendering

**static \_bound(x, min\_value, max\_value)**  
Method used to bound state and action variables.

#### Parameters

- **x** – the variable to bound;
- **min\_value** – the minimum value;
- **max\_value** – the maximum value;

**Returns** The bounded variable.

**info**  
An object containing the info of the environment.

**Type** Returns

**seed(seed)**  
Set the seed of the environment.

**Parameters seed** (*float*) – the value of the seed.

## 3.7.2 Generators

## Grid world

```
mushroom_rl.environments.generators.grid_world.generate_grid_world(grid, prob,
pos_rew,
neg_rew,
gamma=0.9,
hori-
zon=100)
```

This Grid World generator requires a .txt file to specify the shape of the grid world and the cells. There are five types of cells: ‘S’ is the starting position where the agent is; ‘G’ is the goal state; ‘.’ is a normal cell; ‘\*’ is a hole, when the agent steps on a hole, it receives a negative reward and the episode ends; ‘#’ is a wall, when the agent is supposed to step on a wall, it actually remains in its current state. The initial states distribution is uniform among all the initial states provided.

The grid is expected to be rectangular.

### Parameters

- **grid** (*str*) – the path of the file containing the grid structure;
- **prob** (*float*) – probability of success of an action;
- **pos\_rew** (*float*) – reward obtained in goal states;
- **neg\_rew** (*float*) – reward obtained in “hole” states;
- **gamma** (*float*, 9) – discount factor;
- **horizon** (*int*, 100) – the horizon.

**Returns** A FiniteMDP object built with the provided parameters.

```
mushroom_rl.environments.generators.grid_world.parse_grid(grid)
```

Parse the grid file:

**Parameters** **grid** (*str*) – the path of the file containing the grid structure;

**Returns** A list containing the grid structure.

```
mushroom_rl.environments.generators.grid_world.compute_probabilities(grid_map,
cell_list,
prob)
```

Compute the transition probability matrix.

### Parameters

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells;
- **prob** (*float*) – probability of success of an action.

**Returns** The transition probability matrix;

```
mushroom_rl.environments.generators.grid_world.compute_reward(grid_map,
cell_list, pos_rew,
neg_rew)
```

Compute the reward matrix.

### Parameters

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells;
- **pos\_rew** (*float*) – reward obtained in goal states;

- **neg\_rew** (*float*) – reward obtained in “hole” states;

**Returns** The reward matrix.

```
mushroom_rl.environments.generators.grid_world.compute_mu(grid_map, cell_list)
Compute the initial states distribution.
```

#### Parameters

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells.

**Returns** The initial states distribution.

### Simple chain

```
mushroom_rl.environments.generators.simple_chain.generate_simple_chain(state_n,
goal_states,
prob,
rew,
mu=None,
gamma=0.9,
horizon=100)
```

Simple chain generator.

#### Parameters

- **state\_n** (*int*) – number of states;
- **goal\_states** (*list*) – list of goal states;
- **prob** (*float*) – probability of success of an action;
- **rew** (*float*) – reward obtained in goal states;
- **mu** (*np.ndarray*) – initial state probability distribution;
- **gamma** (*float*, 0.9) – discount factor;
- **horizon** (*int*, 100) – the horizon.

**Returns** A FiniteMDP object built with the provided parameters.

```
mushroom_rl.environments.generators.simple_chain.compute_probabilities(state_n,
prob)
```

Compute the transition probability matrix.

#### Parameters

- **state\_n** (*int*) – number of states;
- **prob** (*float*) – probability of success of an action.

**Returns** The transition probability matrix;

```
mushroom_rl.environments.generators.simple_chain.compute_reward(state_n,
goal_states,
rew)
```

Compute the reward matrix.

#### Parameters

- **state\_n** (*int*) – number of states;

- **goal\_states** (*list*) – list of goal states;
- **rew** (*float*) – reward obtained in goal states.

**Returns** The reward matrix.

## Taxi

```
mushroom_rl.environments.generators.taxi.generate_taxi(grid, prob=0.9, rew=(0, 1, 3, 15), gamma=0.99, horizon=np.inf)
```

This Taxi generator requires a .txt file to specify the shape of the grid world and the cells. There are five types of cells: ‘S’ is the starting where the agent is; ‘G’ is the goal state; ‘.’ is a normal cell; ‘F’ is a passenger, when the agent steps on a hole, it picks up it. ‘#’ is a wall, when the agent is supposed to step on a wall, it actually remains in its current state. The initial states distribution is uniform among all the initial states provided. The episode terminates when the agent reaches the goal state. The reward is always 0, except for the goal state where it depends on the number of collected passengers. Each action has a certain probability of success and, if it fails, the agent goes in a perpendicular direction from the supposed one.

The grid is expected to be rectangular.

This problem is inspired from: “Bayesian Q-Learning”. Dearden R. et al.. 1998.

### Parameters

- **grid** (*str*) – the path of the file containing the grid structure;
- **prob** (*float*, 0.9) – probability of success of an action;
- **rew** (*tuple*, (0, 1, 3, 15)) – rewards obtained in goal states;
- **gamma** (*float*, 0.99) – discount factor;
- **horizon** (*int*, np.inf) – the horizon.

**Returns** A FiniteMDP object built with the provided parameters.

```
mushroom_rl.environments.generators.taxi.parse_grid(grid)
```

Parse the grid file:

**Parameters** **grid** (*str*) – the path of the file containing the grid structure.

**Returns** A list containing the grid structure.

```
mushroom_rl.environments.generators.taxi.compute_probabilities(grid_map, cell_list, passenger_list, prob)
```

Compute the transition probability matrix.

### Parameters

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells;
- **passenger\_list** (*list*) – list of passenger cells;
- **prob** (*float*) – probability of success of an action.

**Returns** The transition probability matrix;

```
mushroom_rl.environments.generators.taxi.compute_reward(grid_map, cell_list, passenger_list, rew)
```

Compute the reward matrix.

**Parameters**

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells;
- **passenger\_list** (*list*) – list of passenger cells;
- **rew** (*tuple*) – rewards obtained in goal states.

**Returns** The reward matrix.

```
mushroom_rl.environments.generators.taxi.compute_mu(grid_map, cell_list, passenger_list)
```

Compute the initial states distribution.

**Parameters**

- **grid\_map** (*list*) – list containing the grid structure;
- **cell\_list** (*list*) – list of non-wall cells;
- **passenger\_list** (*list*) – list of passenger cells.

**Returns** The initial states distribution.

## 3.8 Features

The features in MushroomRL are 1-D arrays computed applying a specified function to a raw input, e.g. polynomial features of the state of an MDP. MushroomRL supports three types of features:

- basis functions;
- tensor basis functions;
- tiles.

The tensor basis functions are a PyTorch implementation of the standard basis functions. They are less straightforward than the standard ones, but they are faster to compute as they can exploit parallel computing, e.g. GPU-acceleration and multi-core systems.

All the types of features are exposed by a single factory method `Features` that builds the one requested by the user.

```
mushroom_rl.features.features.Features(basis_list=None, tilings=None, tensor_list=None, n_outputs=None, function=None, device=None)
```

Factory method to build the requested type of features. The types are mutually exclusive.

Possible features are tilings (`tilings`), basis functions (`basis_list`), tensor basis (`tensor_list`), and functional mappings (`n_outputs` and `function`).

The difference between `basis_list` and `tensor_list` is that the former is a list of python classes each one evaluating a single element of the feature vector, while the latter consists in a list of PyTorch modules that can be used to build a PyTorch network. The use of `tensor_list` is a faster way to compute features than `basis_list` and is suggested when the computation of the requested features is slow (see the Gaussian radial basis function implementation as an example). A functional mapping applies a function to the input computing an `n_outputs`-dimensional vector, where the mapping is expressed by `function`. If `function` is not provided, the identity is used.

**Parameters**

- **basis\_list** (*list, None*) – list of basis functions;
- **tilings** (*[object, list], None*) – single object or list of tilings;

- **tensor\_list** (*list, None*) – list of dictionaries containing the instructions to build the requested tensors;
- **n\_outputs** (*int, None*) – dimensionality of the feature mapping;
- **function** (*object, None*) – a callable function to be used as feature mapping. Only needed when using a functional mapping.
- **device** (*int, None*) – where to run the group of tensors. Only needed when using a list of tensors.

**Returns** The class implementing the requested type of features.

`mushroom_rl.features.features.get_action_features(phi_state, action, n_actions)`

Compute an array of size  $\text{len}(\phi\text{phi\_state}) * \text{n\_actions}$  filled with zeros, except for elements from  $\text{len}(\phi\text{phi\_state}) * \text{action}$  to  $\text{len}(\phi\text{phi\_state}) * (\text{action} + 1)$  that are filled with  $\phi\text{phi\_state}$ . This is used to compute state-action features.

#### Parameters

- **phi\_state** (*np.ndarray*) – the feature of the state;
- **action** (*np.ndarray*) – the action whose features have to be computed;
- **n\_actions** (*int*) – the number of actions.

**Returns** The state-action features.

The factory method returns a class that extends the abstract class `FeatureImplementation`.

The documentation for every feature type can be found here:

### 3.8.1 Basis

#### Fourier

`class mushroom_rl.features.basis.fourier.FourierBasis(low, delta, c, dimensions=None)`

Bases: `object`

Class implementing Fourier basis functions. The value of the feature is computed using the formula:

$$\sum \cos \pi(X - m) / \Delta c$$

where X is the input, m is the vector of the minimum input values (for each dimensions), Delta is the vector of maximum

`__init__(low, delta, c, dimensions=None)`

Constructor.

#### Parameters

- **low** (*np.ndarray*) – vector of minimum values of the input variables;
- **delta** (*np.ndarray*) – vector of the maximum difference between two values of the input variables, i.e.  $\text{delta} = \text{high} - \text{low}$ ;
- **c** (*np.ndarray*) – vector of weights for the state variables;
- **dimensions** (*list, None*) – list of the dimensions of the input to be considered by the feature.

`__call__(x)`

Call self as a function.

**static generate**(*low*, *high*, *n*, *dimensions=None*)

Factory method to build a set of fourier basis.

**Parameters**

- **low** (*np.ndarray*) – vector of minimum values of the input variables;
- **high** (*np.ndarray*) – vector of maximum values of the input variables;
- **n** (*int*) – number of harmonics to consider for each state variable
- **dimensions** (*list*, *None*) – list of the dimensions of the input to be considered by the features.

**Returns** The list of the generated fourier basis functions.

**Gaussian RBF****class** mushroom\_rl.features.basis.gaussian\_rbf.**GaussianRBF**(*mean*, *scale*, *dimensions=None*)

Bases: object

Class implementing Gaussian radial basis functions. The value of the feature is computed using the formula:

$$\sum \frac{(X_i - \mu_i)^2}{\sigma_i}$$

where X is the input, mu is the mean vector and sigma is the scale parameter vector.

**\_\_init\_\_**(*mean*, *scale*, *dimensions=None*)

Constructor.

**Parameters**

- **mean** (*np.ndarray*) – the mean vector of the feature;
- **scale** (*np.ndarray*) – the scale vector of the feature;
- **dimensions** (*list*, *None*) – list of the dimensions of the input to be considered by the feature. The number of dimensions must match the dimensionality of mean and scale.

**\_\_call\_\_**(*x*)

Call self as a function.

**static generate**(*n\_centers*, *low*, *high*, *dimensions=None*)

Factory method to build uniformly spaced gaussian radial basis functions with a 25% overlap.

**Parameters**

- **n\_centers** (*list*) – list of the number of radial basis functions to be used for each dimension.
- **low** (*np.ndarray*) – lowest value for each dimension;
- **high** (*np.ndarray*) – highest value for each dimension;
- **dimensions** (*list*, *None*) – list of the dimensions of the input to be considered by the feature. The number of dimensions must match the number of elements in n\_centers and low.

**Returns** The list of the generated radial basis functions.

## Polynomial

```
class mushroom_rl.features.basis.polynomial.PolynomialBasis (dimensions=None,
degrees=None)
```

Bases: object

Class implementing polynomial basis functions. The value of the feature is computed using the formula:

$$\prod X_i^{d_i}$$

where X is the input and d is the vector of the exponents of the polynomial.

**\_\_init\_\_** (*dimensions=None, degrees=None*)

Constructor. If both parameters are None, the constant feature is built.

### Parameters

- **dimensions** (*list, None*) – list of the dimensions of the input to be considered by the feature;
- **degrees** (*list, None*) – list of the degrees of each dimension to be considered by the feature. It must match the number of elements of dimensions.

**\_\_call\_\_** (*x*)

Call self as a function.

**static \_compute\_exponents** (*order, n\_variables*)

Find the exponents of a multivariate polynomial expression of order *order* and *n\_variables* number of variables.

### Parameters

- **order** (*int*) – the maximum order of the polynomial;
- **n\_variables** (*int*) – the number of elements of the input vector.

**Yields** The current exponent of the polynomial.

**static generate** (*max\_degree, input\_size*)

Factory method to build a polynomial of order *max\_degree* based on the first *input\_size* dimensions of the input.

### Parameters

- **max\_degree** (*int*) – maximum degree of the polynomial;
- **input\_size** (*int*) – size of the input.

**Returns** The list of the generated polynomial basis functions.

## 3.8.2 Tensors

### Gaussian tensor

```
class mushroom_rl.features.tensors.gaussian_tensor.PyTorchGaussianRBF (mu,
scale,
dim)
```

Bases: sphinx.ext.autodoc.importer.\_MockObject

Pytorch module to implement a gaussian radial basis function.

**\_\_init\_\_** (*mu, scale, dim*)

Initialize self. See help(type(self)) for accurate signature.

**static generate**(*n\_centers*, *low*, *high*, *dimensions*=None)

Factory method that generates the list of dictionaries to build the tensors representing a set of uniformly spaced Gaussian radial basis functions with a 25% overlap.

**Parameters**

- **n\_centers** (*list*) – list of the number of radial basis functions to be used for each dimension;
- **low** (*np.ndarray*) – lowest value for each dimension;
- **high** (*np.ndarray*) – highest value for each dimension;
- **dimensions** (*list, None*) – list of the dimensions of the input to be considered by the feature. The number of dimensions must match the number of elements in *n\_centers* and *low*.

**Returns** The list of dictionaries as described above.

### 3.8.3 Tiles

**class** mushroom\_rl.features.tiles.tiles.**Tiles**(*x\_range*, *n\_tiles*, *state\_components*=None)

Bases: object

Class implementing rectangular tiling. For each point in the state space, this class can be used to compute the index of the corresponding tile.

**\_\_init\_\_**(*x\_range*, *n\_tiles*, *state\_components*=None)

Constructor.

**Parameters**

- **x\_range** (*list*) – list of two-elements lists specifying the range of each state variable;
- **n\_tiles** (*list*) – list of the number of tiles to be used for each dimension.
- **state\_components** (*list, None*) – list of the dimensions of the input to be considered by the tiling. The number of elements must match the number of elements in *x\_range* and *n\_tiles*.

**\_\_call\_\_**(*x*)

Call self as a function.

**static generate**(*n\_tilings*, *n\_tiles*, *low*, *high*, *uniform*=False)

Factory method to build *n\_tilings* tilings of *n\_tiles* tiles with a range between *low* and *high* for each dimension.

**Parameters**

- **n\_tilings** (*int*) – number of tilings;
- **n\_tiles** (*list*) – number of tiles for each tilings for each dimension;
- **low** (*np.ndarray*) – lowest value for each dimension;
- **high** (*np.ndarray*) – highest value for each dimension.
- **uniform** (*bool, False*) – if True the displacement for each tiling will be w/n\_tilings, where w is the tile width. Otherwise, the displacement will be k\*w/n\_tilings, where k=2i+1, where i is the dimension index.

**Returns** The list of the generated tiles.

## 3.9 Policy

```
class mushroom_rl.policy.policy.Policy
Bases: object
```

Interface representing a generic policy. A policy is a probability distribution that gives the probability of taking an action given a specified state. A policy is used by mushroom agents to interact with the environment.

```
__call__(*args)
```

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args` (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

```
draw_action(state)
```

Sample an action in state using the policy.

**Parameters** `state` (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

```
reset()
```

Useful when the policy needs a special initialization at the beginning of an episode.

```
__init__
```

Initialize self. See help(type(self)) for accurate signature.

```
class mushroom_rl.policy.policy.ParametricPolicy
Bases: mushroom_rl.policy.policy.Policy
```

Interface for a generic parametric policy. A parametric policy is a policy that depends on set of parameters, called the policy weights. If the policy is differentiable, the derivative of the probability for a specified state-action pair can be provided.

```
diff_log(state, action)
```

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (*np.ndarray*) – the state where the gradient is computed
- `action` (*np.ndarray*) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

```
diff(state, action)
```

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (*np.ndarray*) – the state where the derivative is computed
- `action` (*np.ndarray*) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

**set\_weights** (*weights*)  
Setter.

**Parameters** **weights** (*np.ndarray*) – the vector of the new weights to be used by the policy.

**get\_weights** ()  
Getter.

**Returns** The current policy weights.

**weights\_size**  
Property.

**Returns** The size of the policy weights.

**\_\_call\_\_** (\**args*)  
Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**\_\_init\_\_**  
Initialize self. See help(type(self)) for accurate signature.

**draw\_action** (*state*)  
Sample an action in *state* using the policy.

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**reset** ()  
Useful when the policy needs a special initialization at the beginning of an episode.

### 3.9.1 Deterministic policy

```
class mushroom_rl.policy.deterministic_policy.DeterministicPolicy(mu)
Bases: mushroom_rl.policy.ParametricPolicy
```

Simple parametric policy representing a deterministic policy. As deterministic policies are degenerate probability functions where all the probability mass is on the deterministic action, they are not differentiable, even if the mean value approximator is differentiable.

**\_\_init\_\_** (*mu*)  
Constructor.

**Parameters** **mu** (*Regressor*) – the regressor representing the action to select in each state.

**get\_regressor** ()  
Getter.

**Returns** the regressor that is used to map state to actions.

**\_\_call\_\_** (*state, action*)  
Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**draw\_action(state)**

Sample an action in `state` using the policy.

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action sampled from the policy.

**set\_weights(weights)**

Setter.

**Parameters** `weights` (`np.ndarray`) – the vector of the new weights to be used by the policy.

**get\_weights()**

Getter.

**Returns** The current policy weights.

**weights\_size**

Property.

**Returns** The size of the policy weights.

**diff(state, action)**

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (`np.ndarray`) – the state where the derivative is computed
- `action` (`np.ndarray`) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

**diff\_log(state, action)**

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (`np.ndarray`) – the state where the gradient is computed
- `action` (`np.ndarray`) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

**reset()**

Useful when the policy needs a special initialization at the beginning of an episode.

### 3.9.2 Gaussian policy

```
class mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy
Bases: mushroom_rl.policy.policy.ParametricPolicy
```

Abstract class of Gaussian policies.

**`__call__(state, action)`**

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args` (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**`draw_action(state)`**

Sample an action in `state` using the policy.

**Parameters** `state` (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**`__init__`**

Initialize self. See `help(type(self))` for accurate signature.

**`diff(state, action)`**

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (*np.ndarray*) – the state where the derivative is computed
- `action` (*np.ndarray*) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

**`diff_log(state, action)`**

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state` (*np.ndarray*) – the state where the gradient is computed
- `action` (*np.ndarray*) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

**`get_weights()`**

Getter.

**Returns** The current policy weights.

**`reset()`**

Useful when the policy needs a special initialization at the beginning of an episode.

**`set_weights(weights)`**

Setter.

**Parameters** `weights` (*np.ndarray*) – the vector of the new weights to be used by the policy.

**`weights_size`**

Property.

**Returns** The size of the policy weights.

```
class mushroom_rl.policy.gaussian_policy.GaussianPolicy(mu, sigma)
Bases: mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy
```

Gaussian policy. This is a differentiable policy for continuous action spaces. The policy samples an action in every state following a gaussian distribution, where the mean is computed in the state and the covariance matrix is fixed.

**\_\_init\_\_(mu, sigma)**

Constructor.

#### Parameters

- **mu** (`Regressor`) – the regressor representing the mean w.r.t. the state;
- **sigma** (`np.ndarray`) – a square positive definite matrix representing the covariance matrix. The size of this matrix must be n x n, where n is the action dimensionality.

**set\_sigma(sigma)**

Setter.

**Parameters** **sigma** (`np.ndarray`) – the new covariance matrix. Must be a square positive definite matrix.

**diff\_log(state, action)**

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

#### Parameters

- **state** (`np.ndarray`) – the state where the gradient is computed
- **action** (`np.ndarray`) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

**set\_weights(weights)**

Setter.

**Parameters** **weights** (`np.ndarray`) – the vector of the new weights to be used by the policy.

**get\_weights()**

Getter.

**Returns** The current policy weights.

**weights\_size**

Property.

**Returns** The size of the policy weights.

**call(state, action)**

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (`list`) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy. If the action space is continuous, state and action must be provided

**diff(state, action)**

Compute the derivative of the probability density function, in the specified state and action pair. Normally

it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

### Parameters

- **state** (`np.ndarray`) – the state where the derivative is computed
- **action** (`np.ndarray`) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

### `draw_action(state)`

Sample an action in `state` using the policy.

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action sampled from the policy.

### `reset()`

Useful when the policy needs a special initialization at the beginning of an episode.

## `class mushroom_rl.policy.gaussian_policy.DiagonalGaussianPolicy(mu, std)`

Bases: `mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy`

Gaussian policy with learnable standard deviation. The Covariance matrix is constrained to be a diagonal matrix, where the diagonal is the squared standard deviation vector. This is a differentiable policy for continuous action spaces. This policy is similar to the gaussian policy, but the weights includes also the standard deviation.

### `__init__(mu, std)`

Constructor.

### Parameters

- **mu** (`Regressor`) – the regressor representing the mean w.r.t. the state;
- **std** (`np.ndarray`) – a vector of standard deviations. The length of this vector must be equal to the action dimensionality.

### `set_std(std)`

Setter.

**Parameters** `std` (`np.ndarray`) – the new standard deviation. Must be a square positive definite matrix.

### `diff_log(state, action)`

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

### Parameters

- **state** (`np.ndarray`) – the state where the gradient is computed
- **action** (`np.ndarray`) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

### `set_weights(weights)`

Setter.

**Parameters** `weights` (`np.ndarray`) – the vector of the new weights to be used by the policy.

**get\_weights()**

Getter.

**Returns** The current policy weights.

**weights\_size**

Property.

**Returns** The size of the policy weights.

**\_\_call\_\_(state, action)**

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args (list)` – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**diff(state, action)**

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- `state (np.ndarray)` – the state where the derivative is computed
- `action (np.ndarray)` – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

**draw\_action(state)**

Sample an action in state using the policy.

**Parameters** `state (np.ndarray)` – the state where the agent is.

**Returns** The action sampled from the policy.

**reset()**

Useful when the policy needs a special initialization at the beginning of an episode.

```
class mushroom_rl.policy.gaussian_policy.StateStdGaussianPolicy(mu, std,
    eps=1e-06)
Bases: mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy
```

Gaussian policy with learnable standard deviation. The Covariance matrix is constrained to be a diagonal matrix, where the diagonal is the squared standard deviation, which is computed for each state. This is a differentiable policy for continuous action spaces. This policy is similar to the diagonal gaussian policy, but a parametric regressor is used to compute the standard deviation, so the standard deviation depends on the current state.

**\_\_init\_\_(mu, std, eps=1e-06)**

Constructor.

**Parameters**

- `mu (Regressor)` – the regressor representing the mean w.r.t. the state;
- `std (Regressor)` – the regressor representing the standard deviations w.r.t. the state.  
The output dimensionality of the regressor must be equal to the action dimensionality;
- `eps (float, 1e-6)` – A positive constant added to the variance to ensure that is always greater than zero.

**diff\_log**(state, action)

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

**Parameters**

- **state** (*np.ndarray*) – the state where the gradient is computed
- **action** (*np.ndarray*) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

**set\_weights**(weights)

Setter.

**Parameters** **weights** (*np.ndarray*) – the vector of the new weights to be used by the policy.

**get\_weights**()

Getter.

**Returns** The current policy weights.

**weights\_size**

Property.

**Returns** The size of the policy weights.

**\_\_call\_\_**(state, action)

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**diff**(state, action)

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- **state** (*np.ndarray*) – the state where the derivative is computed
- **action** (*np.ndarray*) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights

**draw\_action**(state)

Sample an action in state using the policy.

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**reset**()

Useful when the policy needs a special initialization at the beginning of an episode.

```
class mushroom_rl.policy.gaussian_policy.StateLogStdGaussianPolicy(mu,
                                                                    log_std)
Bases: mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy
```

Gaussian policy with learnable standard deviation. The Covariance matrix is constrained to be a diagonal matrix, the diagonal is computed by an exponential transformation of the logarithm of the standard deviation computed in each state. This is a differentiable policy for continuous action spaces. This policy is similar to the State std gaussian policy, but here the regressor represents the logarithm of the standard deviation.

```
__init__(mu, log_std)
```

Constructor.

#### Parameters

- **mu** (`Regressor`) – the regressor representing the mean w.r.t. the state;
- **log\_std** (`Regressor`) – a regressor representing the logarithm of the variance w.r.t. the state. The output dimensionality of the regressor must be equal to the action dimensionality.

```
diff_log(state, action)
```

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

#### Parameters

- **state** (`np.ndarray`) – the state where the gradient is computed
- **action** (`np.ndarray`) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

```
set_weights(weights)
```

Setter.

**Parameters** **weights** (`np.ndarray`) – the vector of the new weights to be used by the policy.

```
get_weights()
```

Getter.

**Returns** The current policy weights.

```
weights_size
```

Property.

**Returns** The size of the policy weights.

```
__call__(state, action)
```

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (`list`) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

```
diff(state, action)
```

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- **state** (*np.ndarray*) – the state where the derivative is computed
- **action** (*np.ndarray*) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights**draw\_action** (*state*)Sample an action in *state* using the policy.**Parameters** **state** (*np.ndarray*) – the state where the agent is.**Returns** The action sampled from the policy.**reset** ()

Useful when the policy needs a special initialization at the beginning of an episode.

### 3.9.3 Noise policy

```
class mushroom_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy (mu, sigma, theta,
dt, x0=None)
```

Bases: *mushroom\_rl.policy.policy.ParametricPolicy*Ornstein-Uhlenbeck process as implemented in: <https://github.com/openai/baselines/blob/master/baselines/ddpg/noise.py>.

This policy is commonly used in the Deep Deterministic Policy Gradient algorithm.

**\_\_init\_\_** (*mu, sigma, theta, dt, x0=None*)

Constructor.

**Parameters**

- **mu** (*Regressor*) – the regressor representing the mean w.r.t. the state;
- **sigma** (*np.ndarray*) – average magnitude of the random fluctuations per square-root time;
- **theta** (*float*) – rate of mean reversion;
- **dt** (*float*) – time interval;
- **x0** (*np.ndarray, None*) – initial values of noise.

**\_\_call\_\_** (*state, action*)

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided**draw\_action** (*state*)Sample an action in *state* using the policy.**Parameters** **state** (*np.ndarray*) – the state where the agent is.**Returns** The action sampled from the policy.**set\_weights** (*weights*)

Setter.

**Parameters** **weights** (*np.ndarray*) – the vector of the new weights to be used by the policy.

**get\_weights()**

Getter.

**Returns** The current policy weights.**weights\_size**

Property.

**Returns** The size of the policy weights.**reset()**

Useful when the policy needs a special initialization at the beginning of an episode.

**diff(state, action)**

Compute the derivative of the probability density function, in the specified state and action pair. Normally it is computed w.r.t. the derivative of the logarithm of the probability density function, exploiting the likelihood ratio trick, i.e.:

$$\nabla_{\theta} p(s, a) = p(s, a) \nabla_{\theta} \log p(s, a)$$

**Parameters**

- **state** (*np.ndarray*) – the state where the derivative is computed
- **action** (*np.ndarray*) – the action where the derivative is computed

**Returns** The derivative w.r.t. the policy weights**diff\_log(state, action)**

Compute the gradient of the logarithm of the probability density function, in the specified state and action pair, i.e.:

$$\nabla_{\theta} \log p(s, a)$$

**Parameters**

- **state** (*np.ndarray*) – the state where the gradient is computed
- **action** (*np.ndarray*) – the action where the gradient is computed

**Returns** The gradient of the logarithm of the pdf w.r.t. the policy weights

### 3.9.4 TD policy

**class** mushroom\_rl.policy.td\_policy.TDPolicyBases: *mushroom\_rl.policy.policy.Policy***\_\_init\_\_()**

Constructor.

**set\_q(approximator)****Parameters** **approximator** (*object*) – the approximator to use.**get\_q()****Returns** The approximator used by the policy.**\_\_call\_\_(\*args)**

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**draw\_action(state)**

Sample an action in `state` using the policy.

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action sampled from the policy.

**reset()**

Useful when the policy needs a special initialization at the beginning of an episode.

**class mushroom\_rl.policy.td\_policy.EpsGreedy(epsilon)**

Bases: `mushroom_rl.policy.td_policy.TDPolicy`

Epsilon greedy policy.

**\_\_init\_\_(epsilon)**

Constructor.

**Parameters** `epsilon` (`Parameter`) – the exploration coefficient. It indicates the probability of performing a random actions in the current step.

**\_\_call\_\_(\*args)**

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args` (`list`) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**draw\_action(state)**

Sample an action in `state` using the policy.

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action sampled from the policy.

**set\_epsilon(epsilon)**

Setter.

**Parameters**

- `epsilon` (`Parameter`) – the exploration coefficient. It indicates the
- `of performing a random actions in the current step.`  
`(probability)`–

**update(\*idx)**

Update the value of the epsilon parameter at the provided index (e.g. in case of different values of epsilon for each visited state according to the number of visits).

**Parameters** `*idx` (`list`) – index of the parameter to be updated.

**get\_q()**

**Returns** The approximator used by the policy.

**reset()**

Useful when the policy needs a special initialization at the beginning of an episode.

**set\_q(approximator)**

**Parameters** `approximator` (*object*) – the approximator to use.

```
class mushroom_rl.policy.td_policy.Boltzmann(beta)
Bases: mushroom_rl.policy.td_policy.TDPolicy
```

Boltzmann softmax policy.

```
__init__(beta)
```

Constructor.

#### Parameters

- `beta` (`Parameter`) – the inverse of the temperature distribution. As **temperature approaches infinity, the policy becomes more and (the) –**
- **random. As the temperature approaches 0.0, the policy becomes (more) –**
- **and more greedy. (more) –**

```
__call__(*args)
```

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args` (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

```
draw_action(state)
```

Sample an action in `state` using the policy.

**Parameters** `state` (`np.ndarray`) – the state where the agent is.

**Returns** The action sampled from the policy.

```
set_beta(beta)
```

Setter.

**Parameters** `beta` (`Parameter`) – the inverse of the temperature distribution.

```
update(*idx)
```

Update the value of the beta parameter at the provided index (e.g. in case of different values of beta for each visited state according to the number of visits).

**Parameters** `*idx` (*list*) – index of the parameter to be updated.

```
get_q()
```

**Returns** The approximator used by the policy.

```
reset()
```

Useful when the policy needs a special initialization at the beginning of an episode.

```
set_q(approximator)
```

**Parameters** `approximator` (*object*) – the approximator to use.

```
class mushroom_rl.policy.td_policy.Mellowmax(omega, beta_min=-10.0, beta_max=10.0)
Bases: mushroom_rl.policy.td_policy.Boltzmann
```

Mellowmax policy. “An Alternative Softmax Operator for Reinforcement Learning”. Asadi K. and Littman M.L.. 2017.

---

**\_\_init\_\_** (*omega*, *beta\_min*=-10.0, *beta\_max*=10.0)  
Constructor.

#### Parameters

- **omega** ([Parameter](#)) – the omega parameter of the policy from which beta of the Boltzmann policy is computed;
- **beta\_min** (*float*, -10.) – one end of the bracketing interval for minimization with Brent's method;
- **beta\_max** (*float*, 10.) – the other end of the bracketing interval for minimization with Brent's method.

**set\_beta** (*beta*)

Setter.

**Parameters** **beta** ([Parameter](#)) – the inverse of the temperature distribution.

**update** (\**idx*)

Update the value of the beta parameter at the provided index (e.g. in case of different values of beta for each visited state according to the number of visits).

**Parameters** \***idx** (*list*) – index of the parameter to be updated.

**\_\_call\_\_** (\**args*)

Compute the probability of taking action in a certain state following the policy.

**Parameters** \***args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**draw\_action** (*state*)

Sample an action in *state* using the policy.

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**get\_q** ()

**Returns** The approximator used by the policy.

**reset** ()

Useful when the policy needs a special initialization at the beginning of an episode.

**set\_q** (*approximator*)

**Parameters** **approximator** (*object*) – the approximator to use.

### 3.9.5 Torch policy

**class** mushroom\_rl.policy.torch\_policy.**TorchPolicy** (*use\_cuda*)  
Bases: *mushroom\_rl.policy.policy*

Interface for a generic PyTorch policy. A PyTorch policy is a policy implemented as a neural network using PyTorch. Functions ending with '\_t' use tensors as input, and also as output when required.

**\_\_init\_\_** (*use\_cuda*)

Constructor.

**Parameters** **use\_cuda** (*bool*) – whether to use cuda or not.

**\_\_call\_\_(state, action)**

Compute the probability of taking action in a certain state following the policy.

**Parameters** `*args` (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**draw\_action(state)**

Sample an action in `state` using the policy.

**Parameters** `state` (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**distribution(state)**

Compute the policy distribution in the given states.

**Parameters** `state` (*np.ndarray*) – the set of states where the distribution is computed.

**Returns** The torch distribution for the provided states.

**entropy(state=None)**

Compute the entropy of the policy.

**Parameters** `state` (*np.ndarray, None*) – the set of states to consider. If the entropy of the policy can be computed in closed form, then `state` can be `None`.

**Returns** The value of the entropy of the policy.

**draw\_action\_t(state)**

Draw an action given a tensor.

**Parameters** `state` (*torch.Tensor*) – set of states.

**Returns** The tensor of the actions to perform in each state.

**log\_prob\_t(state, action)**

Compute the logarithm of the probability of taking `action` in `state`.

**Parameters**

- `state` (*torch.Tensor*) – set of states.
- `action` (*torch.Tensor*) – set of actions.

**Returns** The tensor of log-probability.

**entropy\_t(state=None)**

Compute the entropy of the policy.

**Parameters** `state` (*torch.Tensor*) – the set of states to consider. If the entropy of the policy can be computed in closed form, then `state` can be `None`.

**Returns** The tensor value of the entropy of the policy.

**distribution\_t(state)**

Compute the policy distribution in the given states.

**Parameters** `state` (*torch.Tensor*) – the set of states where the distribution is computed.

**Returns** The torch distribution for the provided states.

**set\_weights(weights)**

Setter.

**Parameters** `weights` (`np.ndarray`) – the vector of the new weights to be used by the policy.

**get\_weights()**  
Getter.

**Returns** The current policy weights.

**parameters()**  
Returns the trainable policy parameters, as expected by torch optimizers.

**Returns** List of parameters to be optimized.

**reset()**  
Useful when the policy needs a special initialization at the beginning of an episode.

**use\_cuda**  
True if the policy is using cuda\_tensors.

```
class mushroom_rl.policy.torch_policy.GaussianTorchPolicy(network,           in-
                                                               put_shape,          out-
                                                               put_shape, std_0=1.0,
                                                               use_cuda=False,
                                                               **params)
```

Bases: `mushroom_rl.policy.torch_policy.TorchPolicy`

Torch policy implementing a Gaussian policy with trainable standard deviation. The standard deviation is not state-dependent.

**\_\_init\_\_(network, input\_shape, output\_shape, std\_0=1.0, use\_cuda=False, \*\*params)**  
Constructor.

**Parameters**

- **network** (`object`) – the network class used to implement the mean regressor;
- **input\_shape** (`tuple`) – the shape of the state space;
- **output\_shape** (`tuple`) – the shape of the action space;
- **std\_0** (`float, 1.`) – initial standard deviation;
- **params** (`dict`) – parameters used by the network constructor.

**draw\_action\_t(state)**  
Draw an action given a tensor.

**Parameters** `state` (`torch.Tensor`) – set of states.

**Returns** The tensor of the actions to perform in each state.

**log\_prob\_t(state, action)**  
Compute the logarithm of the probability of taking `action` in `state`.

**Parameters**

- **state** (`torch.Tensor`) – set of states.
- **action** (`torch.Tensor`) – set of actions.

**Returns** The tensor of log-probability.

**entropy\_t(state=None)**  
Compute the entropy of the policy.

**Parameters** `state` (`torch.Tensor`) – the set of states to consider. If the entropy of the policy can be computed in closed form, then `state` can be None.

**Returns** The tensor value of the entropy of the policy.

**distribution\_t** (*state*)

Compute the policy distribution in the given states.

**Parameters** **state** (*torch.Tensor*) – the set of states where the distribution is computed.

**Returns** The torch distribution for the provided states.

**set\_weights** (*weights*)

Setter.

**Parameters** **weights** (*np.ndarray*) – the vector of the new weights to be used by the policy.

**get\_weights** ()

Getter.

**Returns** The current policy weights.

**parameters** ()

Returns the trainable policy parameters, as expected by torch optimizers.

**Returns** List of parameters to be optimized.

**\_\_call\_\_** (*state, action*)

Compute the probability of taking action in a certain state following the policy.

**Parameters** **\*args** (*list*) – list containing a state or a state and an action.

**Returns** The probability of all actions following the policy in the given state if the list contains only the state, else the probability of the given action in the given state following the policy.  
If the action space is continuous, state and action must be provided

**distribution** (*state*)

Compute the policy distribution in the given states.

**Parameters** **state** (*np.ndarray*) – the set of states where the distribution is computed.

**Returns** The torch distribution for the provided states.

**draw\_action** (*state*)

Sample an action in *state* using the policy.

**Parameters** **state** (*np.ndarray*) – the state where the agent is.

**Returns** The action sampled from the policy.

**entropy** (*state=None*)

Compute the entropy of the policy.

**Parameters** **state** (*np.ndarray, None*) – the set of states to consider. If the entropy of the policy can be computed in closed form, then *state* can be None.

**Returns** The value of the entropy of the policy.

**reset** ()

Useful when the policy needs a special initialization at the beginning of an episode.

**use\_cuda**

True if the policy is using cuda\_tensors.

## 3.10 Solvers

### 3.10.1 Dynamic programming

```
mushroom_rl.solvers.dynamic_programming.value_iteration(prob, reward, gamma, eps)
```

Value iteration algorithm to solve a dynamic programming problem.

#### Parameters

- **prob** (*np.ndarray*) – transition probability matrix;
- **reward** (*np.ndarray*) – reward matrix;
- **gamma** (*float*) – discount factor;
- **eps** (*float*) – accuracy threshold.

**Returns** The optimal value of each state.

```
mushroom_rl.solvers.dynamic_programming.policy_iteration(prob, reward, gamma)
```

Policy iteration algorithm to solve a dynamic programming problem.

#### Parameters

- **prob** (*np.ndarray*) – transition probability matrix;
- **reward** (*np.ndarray*) – reward matrix;
- **gamma** (*float*) – discount factor.

**Returns** The optimal value of each state and the optimal policy.

### 3.10.2 Car-On-Hill brute-force solver

```
mushroom_rl.solvers.car_on_hill.step(mdp, state, action)
```

Perform a step in the tree.

#### Parameters

- **mdp** ([CarOnHill](#)) – the Car-On-Hill environment;
- **state** (*np.array*) – the state;
- **action** (*np.array*) – the action.

**Returns** The resulting transition executing `action` in `state`.

```
mushroom_rl.solvers.car_on_hill.bfs(mdp, frontier, k, max_k)
```

Perform Breadth-First tree search.

#### Parameters

- **mdp** ([CarOnHill](#)) – the Car-On-Hill environment;
- **frontier** (*list*) – the state at the frontier of the BFS;
- **k** (*int*) – the current depth of the tree;
- **max\_k** (*int*) – maximum depth to consider.

**Returns** A tuple containing a flag for the algorithm ending, and the updated depth of the tree.

```
mushroom_rl.solvers.car_on_hill.solve_car_on_hill(mdp, states, actions, gamma,  
max_k=50)
```

Solver of the Car-On-Hill environment.

#### Parameters

- **mdp** ([CarOnHill](#)) – the Car-On-Hill environment;
- **states** (*np.ndarray*) – the states;
- **actions** (*np.ndarray*) – the actions;
- **gamma** (*float*) – the discount factor;
- **max\_k** (*int*, 50) – maximum depth to consider.

**Returns** The Q-value for each state-action tuple.

## 3.11 Utils

### 3.11.1 Angles

```
mushroom_rl.utils.angles.normalize_angle_positive(angle)
```

Wrap the angle between 0 and  $2 * \pi$ .

**Parameters** **angle** (*float*) – angle to wrap.

**Returns** The wrapped angle.

```
mushroom_rl.utils.angles.normalize_angle(angle)
```

Wrap the angle between  $-\pi$  and  $\pi$ .

**Parameters** **angle** (*float*) – angle to wrap.

**Returns** The wrapped angle.

```
mushroom_rl.utils.angles.shortest_angular_distance(from_angle, to_angle)
```

Compute the shortest distance between two angles

#### Parameters

- **from\_angle** (*float*) – starting angle;
- **to\_angle** (*float*) – final angle.

**Returns** The shortest distance between from\_angle and to\_angle.

```
mushroom_rl.utils.angles.quat_to_euler(quat)
```

Convert a quaternion to euler angles.

**Parameters** **quat** (*np.ndarray*) – quaternion to be converted, must be in format [w, x, y, z]

**Returns** The euler angles [x, y, z] representation of the quaternion

```
mushroom_rl.utils.angles.euler_to_quat(euler)
```

Convert euler angles into a quaternion.

**Parameters** **euler** (*np.ndarray*) – euler angles to be converted

**Returns** Quaternion in format [w, x, y, z]

### 3.11.2 Callbacks

```
class mushroom_rl.utils.callbacks.Callback
Bases: object
```

Interface for all basic callbacks. Implements a list in which it is possible to store data and methods to query and clean the content stored by the callback.

**\_\_init\_\_()**

Constructor.

**\_\_call\_\_(dataset)**

Add samples to the samples list.

**Parameters** **dataset** (*list*) – the samples to collect.

**get()**

**Returns** The current collected data as a list.

**clean()**

Delete the current stored data list

```
class mushroom_rl.utils.callbacks.CollectDataset
```

Bases: mushroom\_rl.utils.callbacks.callback.Callback

This callback can be used to collect samples during the learning of the agent.

**\_\_call\_\_(dataset)**

Add samples to the samples list.

**Parameters** **dataset** (*list*) – the samples to collect.

```
class mushroom_rl.utils.callbacks.CollectQ(approximator)
```

Bases: mushroom\_rl.utils.callbacks.callback.Callback

This callback can be used to collect the action values in all states at the current time step.

**\_\_init\_\_(approximator)**

Constructor.

**Parameters** **approximator** ([[Table](#), [EnsembleTable](#)]) – the approximator to use to predict the action values.

**\_\_call\_\_(dataset)**

Add samples to the samples list.

**Parameters** **dataset** (*list*) – the samples to collect.

```
class mushroom_rl.utils.callbacks.CollectMaxQ(approximator, state)
```

Bases: mushroom\_rl.utils.callbacks.callback.Callback

This callback can be used to collect the maximum action value in a given state at each call.

**\_\_init\_\_(approximator, state)**

Constructor.

**Parameters**

- **approximator** ([[Table](#), [EnsembleTable](#)]) – the approximator to use;
- **state** (*np.ndarray*) – the state to consider.

**\_\_call\_\_(dataset)**

Add samples to the samples list.

**Parameters** `dataset` (*list*) – the samples to collect.

```
class mushroom_rl.utils.callbacks.CollectParameters(parameter, *idx)
Bases: mushroom_rl.utils.callbacks.callback.Callback
```

This callback can be used to collect the values of a parameter (e.g. learning rate) during a run of the agent.

```
__init__(parameter, *idx)
```

Constructor.

#### Parameters

- `parameter` (`Parameter`) – the parameter whose values have to be collected;
- `*idx` (*list*) – index of the parameter when the `parameter` is tabular.

```
__call__(dataset)
```

Add samples to the samples list.

**Parameters** `dataset` (*list*) – the samples to collect.

### 3.11.3 Dataset

```
mushroom_rl.utils.dataset.parse_dataset(dataset, features=None)
```

Split the dataset in its different components and return them.

#### Parameters

- `dataset` (*list*) – the dataset to parse;
- `features` (*object, None*) – features to apply to the states.

**Returns** The `np.ndarray` of state, action, reward, next\_state, absorbing flag and last step flag. Features are applied to state and next\_state, when provided.

```
mushroom_rl.utils.dataset.arrays_as_dataset(states, actions, rewards, next_states, absorbings, lasts)
```

Creates a dataset of transitions from the provided arrays.

#### Parameters

- `states` (`np.ndarray`) – array of states;
- `actions` (`np.ndarray`) – array of actions;
- `rewards` (`np.ndarray`) – array of rewards;
- `next_states` (`np.ndarray`) – array of next\_states;
- `absorbings` (`np.ndarray`) – array of absorbing flags;
- `lasts` (`np.ndarray`) – array of last flags.

**Returns** The list of transitions.

```
mushroom_rl.utils.dataset.episodes_length(dataset)
```

Compute the length of each episode in the dataset.

**Parameters** `dataset` (*list*) – the dataset to consider.

**Returns** A list of length of each episode in the dataset.

```
mushroom_rl.utils.dataset.select_first_episodes(dataset, n_episodes, parse=False)
```

Return the first `n_episodes` episodes in the provided dataset.

#### Parameters

- **dataset** (*list*) – the dataset to consider;
- **n\_episodes** (*int*) – the number of episodes to pick from the dataset;
- **parse** (*bool, False*) – whether to parse the dataset to return.

**Returns** A subset of the dataset containing the first *n\_episodes* episodes.

`mushroom_rl.utils.dataset.select_random_samples(dataset, n_samples, parse=False)`

Return the randomly picked desired number of samples in the provided dataset.

#### Parameters

- **dataset** (*list*) – the dataset to consider;
- **n\_samples** (*int*) – the number of samples to pick from the dataset;
- **parse** (*bool, False*) – whether to parse the dataset to return.

**Returns** A subset of the dataset containing randomly picked *n\_samples* samples.

`mushroom_rl.utils.dataset.compute_J(dataset, gamma=1.0)`

Compute the cumulative discounted reward of each episode in the dataset.

#### Parameters

- **dataset** (*list*) – the dataset to consider;
- **gamma** (*float, 1.*) – discount factor.

**Returns** The cumulative discounted reward of each episode in the dataset.

`mushroom_rl.utils.dataset.compute_metrics(dataset, gamma=1.0)`

Compute the metrics of each complete episode in the dataset.

#### Parameters

- **dataset** (*list*) – the dataset to consider;
- **gamma** (*float, 1.*) – the discount factor.

#### Returns

The minimum score reached in an episode, the maximum score reached in an episode, the mean score reached, the number of completed games.

If episode has not been completed, it returns 0 for all values.

### 3.11.4 Eligibility trace

`mushroom_rl.utils.eligibility_trace.EligibilityTrace(shape, name='replacing')`

Factory method to create an eligibility trace of the provided type.

#### Parameters

- **shape** (*list*) – shape of the eligibility trace table;
- **name** (*str, 'replacing'*) – type of the eligibility trace.

**Returns** The eligibility trace table of the provided shape and type.

`class mushroom_rl.utils.eligibility_trace.ReplacingTrace(shape, initial_value=0.0, dtype=None)`

Bases: `mushroom_rl.utils.table.Table`

Replacing trace.

**reset()**

**update(state, action)**

**\_\_init\_\_(shape, initial\_value=0.0, dtype=None)**  
Constructor.

#### Parameters

- **shape** (*tuple*) – the shape of the tabular regressor.
- **initial\_value** (*float*, *0.*) – the initial value for each entry of the tabular regressor.
- **dtype** (*[int, float]*, *None*) – the dtype of the table array.

**fit(x, y)**

#### Parameters

- **x** (*int*) – index of the table to be filled;
- **y** (*float*) – value to fill in the table.

**n\_actions**

The number of actions considered by the table.

#### Type Returns

**predict(\*z)**

Predict the output of the table given an input.

#### Parameters

- **\*z** (*list*) – list of input of the model. If the table is a Q-table,
- **list may contain states or states and actions depending (this)** – on whether the call requires to predict all q-values or only one q-value corresponding to the provided action;

#### Returns

The table prediction.

**shape**

The shape of the table.

#### Type Returns

**class mushroom\_rl.utils.eligibility\_trace.AccumulatingTrace(shape, initial\_value=0.0, dtype=None)**

Bases: *mushroom\_rl.utils.table.Table*

Accumulating trace.

**reset()**

**update(state, action)**

**\_\_init\_\_(shape, initial\_value=0.0, dtype=None)**  
Constructor.

#### Parameters

- **shape** (*tuple*) – the shape of the tabular regressor.
- **initial\_value** (*float*, *0.*) – the initial value for each entry of the tabular regressor.

- **dtype** (*[int, float], None*) – the dtype of the table array.

**fit** (*x, y*)

#### Parameters

- **x** (*int*) – index of the table to be filled;
- **y** (*float*) – value to fill in the table.

**n\_actions**

The number of actions considered by the table.

#### Type Returns

**predict** (\**z*)

Predict the output of the table given an input.

#### Parameters

- **\*z** (*list*) – list of input of the model. If the table is a Q-table,
- **list may contain states or states and actions depending (this)** – on whether the call requires to predict all q-values or only one q-value corresponding to the provided action;

**Returns** The table prediction.

**shape**

The shape of the table.

#### Type Returns

### 3.11.5 Features

`mushroom_rl.utils.features.uniform_grid(n_centers, low, high)`

This function is used to create the parameters of uniformly spaced radial basis functions with 25% of overlap. It creates a uniformly spaced grid of `n_centers[i]` points in each `ranges[i]`. Also returns a vector containing the appropriate scales of the radial basis functions.

#### Parameters

- **n\_centers** (*list*) – number of centers of each dimension;
- **low** (*np.ndarray*) – lowest value for each dimension;
- **high** (*np.ndarray*) – highest value for each dimension.

**Returns** The uniformly spaced grid and the scale vector.

### 3.11.6 Folder

`mushroom_rl.utils.folder.mk_dir_recursive(dir_path)`

Create a directory and, if needed, all the directory tree. Differently from `os.mkdir`, this function does not raise exception when the directory already exists.

**Parameters** `dir_path` (*str*) – the path of the directory to create.

`mushroom_rl.utils.folder.force_symlink(src, dst)`

Create a symlink deleting the previous one, if it already exists.

#### Parameters

- **src** (*str*) – source;
- **dst** (*str*) – destination.

### 3.11.7 Minibatches

`mushroom_rl.utils.minibatches.minibatch_number(size, batch_size)`

Function to retrieve the number of batches, given a batch sizes.

#### Parameters

- **size** (*int*) – size of the dataset;
- **batch\_size** (*int*) – size of the batches.

**Returns** The number of minibatches in the dataset.

`mushroom_rl.utils.minibatches.minibatch_generator(batch_size, *dataset)`

Generator that creates a minibatch from the full dataset.

#### Parameters

- **batch\_size** (*int*) – the maximum size of each minibatch;
- **dataset** – the dataset to be splitted.

**Returns** The current minibatch.

### 3.11.8 Numerical gradient

`mushroom_rl.utils.numerical_gradient.numerical_diff_policy(policy, state, action, eps=1e-06)`

Compute the gradient of a policy in (state, action) numerically.

#### Parameters

- **policy** (`Policy`) – the policy whose gradient has to be returned;
- **state** (`np.ndarray`) – the state;
- **action** (`np.ndarray`) – the action;
- **eps** (`float, 1e-6`) – the value of the perturbation.

**Returns** The gradient of the provided policy in (state, action) computed numerically.

`mushroom_rl.utils.numerical_gradient.numerical_diff_dist(dist, theta, eps=1e-06)`

Compute the gradient of a distribution in theta numerically.

#### Parameters

- **dist** (`Distribution`) – the distribution whose gradient has to be returned;
- **theta** (`np.ndarray`) – the parametrization where to compute the gradient;
- **eps** (`float, 1e-6`) – the value of the perturbation.

**Returns** The gradient of the provided distribution theta computed numerically.

### 3.11.9 Parameters

```
class mushroom_rl.utils.parameters.Parameter(value, min_value=None, max_value=None,  
                                              size=(1, ))
```

Bases: object

This class implements function to manage parameters, such as learning rate. It also allows to have a single parameter for each state of state-action tuple.

```
__init__(value, min_value=None, max_value=None, size=(1, ))  
Constructor.
```

#### Parameters

- **value** (*float*) – initial value of the parameter;
- **min\_value** (*float*, *None*) – minimum value that the parameter can reach when decreasing;
- **max\_value** (*float*, *None*) – maximum value that the parameter can reach when increasing;
- **size** (*tuple*, *(1, )*) – shape of the matrix of parameters; this shape can be used to have a single parameter for each state or state-action tuple.

```
_call__(*idx, **kwargs)
```

Update and return the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

```
get_value(*idx, **kwargs)
```

Return the current value of the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

```
_compute(*idx, **kwargs)
```

**Returns** The value of the parameter in the provided index.

```
update(*idx, **kwargs)
```

Updates the number of visit of the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter whose number of visits has to be updated.

#### shape

The shape of the table of parameters.

**Type** Returns

```
class mushroom_rl.utils.parameters.LinearParameter(value, threshold_value, n, size=(1,  
                                              ))
```

Bases: *mushroom\_rl.utils.parameters.Parameter*

This class implements a linearly changing parameter according to the number of times it has been used.

```
__init__(value, threshold_value, n, size=(1, ))  
Constructor.
```

#### Parameters

- **value** (*float*) – initial value of the parameter;

- **min\_value** (*float, None*) – minimum value that the parameter can reach when decreasing;
- **max\_value** (*float, None*) – maximum value that the parameter can reach when increasing;
- **size** (*tuple, (1,)*) – shape of the matrix of parameters; this shape can be used to have a single parameter for each state or state-action tuple.

#### `_compute(*idx, **kwargs)`

Returns: The value of the parameter in the provided index.

#### `__call__(*idx, **kwargs)`

Update and return the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

#### `get_value(*idx, **kwargs)`

Return the current value of the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

#### `shape`

The shape of the table of parameters.

**Type** Returns

#### `update(*idx, **kwargs)`

Updates the number of visit of the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter whose number of visits has to be updated.

```
class mushroom_rl.utils.parameters.ExponentialParameter(value, exp=1.0,
                                                       min_value=None,
                                                       max_value=None, size=(1,))
```

Bases: *mushroom\_rl.utils.parameters.Parameter*

This class implements a exponentially changing parameter according to the number of times it has been used.

#### `__init__(value, exp=1.0, min_value=None, max_value=None, size=(1,))`

Constructor.

#### Parameters

- **value** (*float*) – initial value of the parameter;
- **min\_value** (*float, None*) – minimum value that the parameter can reach when decreasing;
- **max\_value** (*float, None*) – maximum value that the parameter can reach when increasing;
- **size** (*tuple, (1,)*) – shape of the matrix of parameters; this shape can be used to have a single parameter for each state or state-action tuple.

#### `_compute(*idx, **kwargs)`

Returns: The value of the parameter in the provided index.

#### `__call__(*idx, **kwargs)`

Update and return the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

**get\_value** (\**idx*, \*\**kwargs*)

Return the current value of the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

**shape**

The shape of the table of parameters.

**Type** Returns

**update** (\**idx*, \*\**kwargs*)

Updates the number of visit of the parameter in the provided index.

**Parameters** `*idx` (*list*) – index of the parameter whose number of visits has to be updated.

**class** mushroom\_rl.utils.parameters.**AdaptiveParameter** (*value*)

Bases: object

This class implements a basic adaptive gradient step. Instead of moving of a step proportional to the gradient, takes a step limited by a given metric. To specify the metric, the natural gradient has to be provided. If natural gradient is not provided, the identity matrix is used.

The step rule is:

$$\begin{aligned} \Delta\theta = \underset{\Delta\vartheta}{\operatorname{argmax}} \Delta\vartheta^t \nabla_{\theta} J \\ \text{s.t. : } \Delta\vartheta^T M \Delta\vartheta \leq \varepsilon \end{aligned}$$

Lecture notes, Neumann G. <http://www.ias.informatik.tu-darmstadt.de/uploads/Geri/lecture-notes-constraint.pdf>

**\_\_init\_\_** (*value*)

Initialize self. See help(type(self)) for accurate signature.

**\_\_call\_\_** (\**args*, \*\**kwargs*)

Call self as a function.

### 3.11.10 Replay memory

**class** mushroom\_rl.utils.replay\_memory.**ReplayMemory** (*initial\_size*, *max\_size*)

Bases: object

This class implements function to manage a replay memory as the one used in “Human-Level Control Through Deep Reinforcement Learning” by Mnih V. et al..

**\_\_init\_\_** (*initial\_size*, *max\_size*)

Constructor.

**Parameters**

- **initial\_size** (*int*) – initial number of elements in the replay memory;

- **max\_size** (*int*) – maximum number of elements that the replay memory can contain.

**add** (*dataset*)

Add elements to the replay memory.

**Parameters** `dataset` (*list*) – list of elements to add to the replay memory.

**get** (*n\_samples*)

Returns the provided number of states from the replay memory.

**Parameters** **n\_samples** (*int*) – the number of samples to return.

**Returns** The requested number of samples.

**reset** ()

Reset the replay memory.

**initialized**

Whether the replay memory has reached the number of elements that allows it to be used.

**Type** Returns

**size**

The number of elements contained in the replay memory.

**Type** Returns

**class** mushroom\_rl.utils.replay\_memory.**SumTree** (*max\_size*)

Bases: object

This class implements a sum tree data structure. This is used, for instance, by PrioritizedReplayMemory.

**\_\_init\_\_** (*max\_size*)

Constructor.

**Parameters** **max\_size** (*int*) – maximum size of the tree.

**add** (*dataset, priority*)

Add elements to the tree.

**Parameters**

- **dataset** (*list*) – list of elements to add to the tree;
- **p** (*np.ndarray*) – priority of each sample in the dataset.

**get** (*s*)

Returns the provided number of states from the replay memory.

**Parameters** **s** (*float*) – the value of the samples to return.

**Returns** The requested sample.

**update** (*idx, priorities*)

Update the priority of the sample at the provided index in the dataset.

**Parameters**

- **idx** (*np.ndarray*) – indexes of the transitions in the dataset;
- **priorities** (*np.ndarray*) – priorities of the transitions.

**size**

The current size of the tree.

**Type** Returns

**max\_p**

The maximum priority among the ones in the tree.

**Type** Returns

**total\_p**

The sum of the priorities in the tree, i.e. the value of the root node.

**Type** Returns

```
class mushroom_rl.utils.replay_memory.PrioritizedReplayMemory(initial_size,
                                                               max_size,      al-
                                                               pha,   beta,   ep-
                                                               silon=0.01)
```

Bases: object

This class implements function to manage a prioritized replay memory as the one used in “Prioritized Experience Replay” by Schaul et al., 2015.

**\_\_init\_\_** (initial\_size, max\_size, alpha, beta, epsilon=0.01)

Constructor.

#### Parameters

- **initial\_size** (*int*) – initial number of elements in the replay memory;
- **max\_size** (*int*) – maximum number of elements that the replay memory can contain;
- **alpha** (*float*) – prioritization coefficient;
- **beta** (*float*) – importance sampling coefficient;
- **epsilon** (*float*, 0.01) – small value to avoid zero probabilities.

**add** (dataset, p)

Add elements to the replay memory.

#### Parameters

- **dataset** (*list*) – list of elements to add to the replay memory;
- **p** (*np.ndarray*) – priority of each sample in the dataset.

**get** (n\_samples)

Returns the provided number of states from the replay memory.

**Parameters** **n\_samples** (*int*) – the number of samples to return.

**Returns** The requested number of samples.

**update** (error, idx)

Update the priority of the sample at the provided index in the dataset.

#### Parameters

- **error** (*np.ndarray*) – errors to consider to compute the priorities;
- **idx** (*np.ndarray*) – indexes of the transitions in the dataset.

#### initialized

Whether the replay memory has reached the number of elements that allows it to be used.

**Type** Returns

#### max\_priority

The maximum value of priority inside the replay memory.

**Type** Returns

### 3.11.11 Spaces

```
class mushroom_rl.utils.spaces.Box(low, high, shape=None)
Bases: object
```

This class implements functions to manage continuous states and action spaces. It is similar to the `Box` class in `gym.spaces.box`.

```
__init__(low, high, shape=None)
Constructor.
```

#### Parameters

- **low** (`[float, np.ndarray]`) – the minimum value of each dimension of the space. If a scalar value is provided, this value is considered as the minimum one for each dimension. If a `np.ndarray` is provided, each i-th element is considered the minimum value of the i-th dimension;
- **high** (`[float, np.ndarray]`) – the maximum value of dimensions of the space. If a scalar value is provided, this value is considered as the maximum one for each dimension. If a `np.ndarray` is provided, each i-th element is considered the maximum value of the i-th dimension;
- **shape** (`np.ndarray, None`) – the dimension of the space. Must match the shape of `low` and `high`, if they are `np.ndarray`.

#### low

The minimum value of each dimension of the space.

**Type** Returns

#### high

The maximum value of each dimension of the space.

**Type** Returns

#### shape

The dimensions of the space.

**Type** Returns

```
class mushroom_rl.utils.spaces.Discrete(n)
Bases: object
```

This class implements functions to manage discrete states and action spaces. It is similar to the `Discrete` class in `gym.spaces.discrete`.

```
__init__(n)
Constructor.
```

**Parameters** `n` (`int`) – the number of values of the space.

#### size

The number of elements of the space.

**Type** Returns

#### shape

The shape of the space that is always `(1,)`.

**Type** Returns

### 3.11.12 Table

**class** mushroom\_rl.utils.table.Table(*shape*, *initial\_value*=0.0, *dtype*=None)  
Bases: object

Table regressor. Used for discrete state and action spaces.

**\_\_init\_\_**(*shape*, *initial\_value*=0.0, *dtype*=None)  
Constructor.

#### Parameters

- **shape** (*tuple*) – the shape of the tabular regressor.
- **initial\_value** (*float*, 0.) – the initial value for each entry of the tabular regressor.
- **dtype** ([*int*, *float*], *None*) – the dtype of the table array.

**fit**(*x*, *y*)

#### Parameters

- **x** (*int*) – index of the table to be filled;
- **y** (*float*) – value to fill in the table.

**predict**(\**z*)

Predict the output of the table given an input.

#### Parameters

- **\*z** (*list*) – list of input of the model. If the table is a Q-table,
- **list may contain states or states and actions depending (this)** – on whether the call requires to predict all q-values or only one q-value corresponding to the provided action;

**Returns** The table prediction.

**n\_actions**

The number of actions considered by the table.

**Type** Returns

**shape**

The shape of the table.

**Type** Returns

**class** mushroom\_rl.utils.table.EnsembleTable(*n\_models*, *shape*)

Bases: mushroom\_rl.approximators.\_implementations.ensemble.Ensemble

This class implements functions to manage table ensembles.

**\_\_init\_\_**(*n\_models*, *shape*)  
Constructor.

#### Parameters

- **n\_models** (*int*) – number of models in the ensemble;
- **shape** (*np.ndarray*) – shape of each table in the ensemble.

**fit**(\**z*, *idx*=None, \*\**fit\_params*)

Fit the *idx*-th model of the ensemble if *idx* is provided, every model otherwise.

### Parameters

- **\*z** (*list*) – a list containing the inputs to use to predict with each regressor of the ensemble;
- **idx** (*int, None*) – index of the model to fit;
- **\*\*fit\_params** (*dict*) – other params.

### model

The list of the models in the ensemble.

#### Type Returns

**predict** (\*z, idx=None, prediction='mean', compute\_variance=False, \*\*predict\_params)  
Predict.

### Parameters

- **\*z** (*list*) – a list containing the inputs to use to predict with each regressor of the ensemble;
- **idx** (*int, None*) – index of the model to use for prediction;
- **prediction** (*str, 'mean'*) – the type of prediction to make. It can be a ‘mean’ of the ensembles, or a ‘sum’;
- **compute\_variance** (*bool, False*) – whether to compute the variance of the prediction or not;
- **\*\*predict\_params** (*dict*) – other parameters used by the predict method the regressor.

**Returns** The predictions of the model.

### reset()

Reset the model parameters.

## 3.11.13 Torch

mushroom\_rl.utils.torch.set\_weights (*parameters, weights, use\_cuda*)

Function used to set the value of a set of torch parameters given a vector of values.

### Parameters

- **parameters** (*list*) – list of parameters to be considered;
- **weights** (*numpy.ndarray*) – array of the new values for the parameters;
- **use\_cuda** (*bool*) – whether the parameters are cuda tensors or not;

mushroom\_rl.utils.torch.get\_weights (*parameters*)

Function used to get the value of a set of torch parameters as a single vector of values.

**Parameters** **parameters** (*list*) – list of parameters to be considered.

**Returns** A numpy vector consisting of all the values of the vectors.

mushroom\_rl.utils.torch.zero\_grad (*parameters*)

Function used to set to zero the value of the gradient of a set of torch parameters.

**Parameters** **parameters** (*list*) – list of parameters to be considered.

mushroom\_rl.utils.torch.get\_gradient (*params*)

Function used to get the value of the gradient of a set of torch parameters.

**Parameters** `parameters` (*list*) – list of parameters to be considered.

```
mushroom_rl.utils.torch.to_float_tensor(x, use_cuda=False)
```

Function used to convert a numpy array to a float torch tensor.

#### Parameters

- `x` (*np.ndarray*) – numpy array to be converted as torch tensor;
- `use_cuda` (*bool*) – whether to build a cuda tensors or not.

**Returns** A float tensor build from the values contained in the input array.

### 3.11.14 Value Functions

```
mushroom_rl.utils.value_functions.compute_advantage_montecarlo(V, s, ss, r,
                                                               absorbing,
                                                               gamma)
```

Function to estimate the advantage and new value function target over a dataset. The value function is estimated using rollouts (monte carlo estimation).

#### Parameters

- `V` ([Regressor](#)) – the current value function regressor;
- `s` (*numpy.ndarray*) – the set of states in which we want to evaluate the advantage;
- `ss` (*numpy.ndarray*) – the set of next states in which we want to evaluate the advantage;
- `r` (*numpy.ndarray*) – the reward obtained in each transition from state `s` to state `ss`;
- `absorbing` (*numpy.ndarray*) – an array of boolean flags indicating if the reached state is absorbing;
- `gamma` (*float*) – the discount factor of the considered problem.

**Returns** The new estimate for the value function of the next state and the advantage function.

```
mushroom_rl.utils.value_functions.compute_advantage(V, s, ss, r, absorbing, gamma)
```

Function to estimate the advantage and new value function target over a dataset. The value function is estimated using bootstrapping.

#### Parameters

- `V` ([Regressor](#)) – the current value function regressor;
- `s` (*numpy.ndarray*) – the set of states in which we want to evaluate the advantage;
- `ss` (*numpy.ndarray*) – the set of next states in which we want to evaluate the advantage;
- `r` (*numpy.ndarray*) – the reward obtained in each transition from state `s` to state `ss`;
- `absorbing` (*numpy.ndarray*) – an array of boolean flags indicating if the reached state is absorbing;
- `gamma` (*float*) – the discount factor of the considered problem.

**Returns** The new estimate for the value function of the next state and the advantage function.

```
mushroom_rl.utils.value_functions.compute_gae(V, s, ss, r, absorbing, last, gamma, lam)
```

Function to compute Generalized Advantage Estimation (GAE) and new value function target over a dataset.

“High-Dimensional Continuous Control Using Generalized Advantage Estimation”. Schulman J. et al.. 2016.

#### Parameters

- **v** (`Regressor`) – the current value function regressor;
- **s** (`numpy.ndarray`) – the set of states in which we want to evaluate the advantage;
- **ss** (`numpy.ndarray`) – the set of next states in which we want to evaluate the advantage;
- **r** (`numpy.ndarray`) – the reward obtained in each transition from state s to state ss;
- **absorbing** (`numpy.ndarray`) – an array of boolean flags indicating if the reached state is absorbing;
- **last** (`numpy.ndarray`) – an array of boolean flags indicating if the reached state is the last of the trajectory;
- **gamma** (`float`) – the discount factor of the considered problem;
- **lam** (`float`) – the value for the lambda coefficient used by GEA algorithm.

**Returns** The new estimate for the value function of the next state and the estimated generalized advantage.

### 3.11.15 Variance parameters

```
class mushroom_rl.utils.variance_parameters.VarianceParameter(value,      expo-
                                                                nential=False,
                                                                min_value=None,
                                                                tol=1.0, size=(1,
                                                                ))
```

Bases: `mushroom_rl.utils.parameters.Parameter`

Abstract class to implement variance-dependent parameters. A target parameter is expected.

**\_\_init\_\_** (`value, exponential=False, min_value=None, tol=1.0, size=(1, )`)  
Constructor.

**Parameters tol** (`float`) – value of the variance of the target variable such that The parameter value is 0.5.

**\_compute** (\*`idx`, \*\*`kwargs`)

Returns: The value of the parameter in the provided index.

**update** (\*`idx`, \*\*`kwargs`)

Updates the value of the parameter in the provided index.

#### Parameters

- **\*idx** (`list`) – index of the parameter whose number of visits has to be updated.
- **target** (`float`) – Value of the target variable;
- **factor** (`float`) – Multiplicative factor for the parameter value, useful when the parameter depend on another parameter value.

**\_\_call\_\_** (\*`idx`, \*\*`kwargs`)

Update and return the parameter in the provided index.

**Parameters \*idx** (`list`) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

**get\_value** (\*`idx`, \*\*`kwargs`)

Return the current value of the parameter in the provided index.

**Parameters \*idx** (`list`) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

### shape

The shape of the table of parameters.

**Type** Returns

```
class mushroom_rl.utils.variance_parameters.VarianceIncreasingParameter(value,
    exponential=False, min_value=None, tol=1.0, size=(1,))
```

Bases: *mushroom\_rl.utils.variance\_parameters.VarianceParameter*

Class implementing a parameter that increases with the target variance.

### \_\_call\_\_ (\*idx, \*\*kwargs)

Update and return the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

### \_\_init\_\_ (value, exponential=False, min\_value=None, tol=1.0, size=(1,))

Constructor.

**Parameters** **tol** (*float*) – value of the variance of the target variable such that The parameter value is 0.5.

### \_compute (\*idx, \*\*kwargs)

Returns: The value of the parameter in the provided index.

### get\_value (\*idx, \*\*kwargs)

Return the current value of the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

### shape

The shape of the table of parameters.

**Type** Returns

### update (\*idx, \*\*kwargs)

Updates the value of the parameter in the provided index.

#### Parameters

- **\*idx** (*list*) – index of the parameter whose number of visits has to be updated.
- **target** (*float*) – Value of the target variable;
- **factor** (*float*) – Multiplicative factor for the parameter value, useful when the parameter depend on another parameter value.

```
class mushroom_rl.utils.variance_parameters.VarianceDecreasingParameter(value,
    exponential=False, min_value=None, tol=1.0, size=(1,))
```

Bases: *mushroom\_rl.utils.variance\_parameters.VarianceParameter*

Class implementing a parameter that decreases with the target variance.

**\_\_call\_\_(idx, \*\*kwargs)**

Update and return the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

**\_\_init\_\_(value, exponential=False, min\_value=None, tol=1.0, size=(1,))**

Constructor.

**Parameters** **tol** (*float*) – value of the variance of the target variable such that The parameter value is 0.5.

**\_compute(idx, \*\*kwargs)**

Returns: The value of the parameter in the provided index.

**get\_value(idx, \*\*kwargs)**

Return the current value of the parameter in the provided index.

**Parameters** **\*idx** (*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

**shape**

The shape of the table of parameters.

**Type** Returns

**update(idx, \*\*kwargs)**

Updates the value of the parameter in the provided index.

**Parameters**

- **\*idx** (*list*) – index of the parameter whose number of visits has to be updated.
- **target** (*float*) – Value of the target variable;
- **factor** (*float*) – Multiplicative factor for the parameter value, useful when the parameter depend on another parameter value.

```
class mushroom_rl.utils.variance_parameters.WindowedVarianceParameter(value,
    exponential=False, min_value=None, tol=1.0, window_size=100, size=(1,))
```

Bases: `mushroom_rl.utils.parameters.Parameter`

Abstract class to implement variance-dependent parameters. A target parameter is expected. differently from the “Variance Parameter” class the variance is computed in a window interval.

**`__init__(value, exponential=False, min_value=None, tol=1.0, window=100, size=(1,))`**  
Constructor.

#### Parameters

- **`tol`** (`float`) – value of the variance of the target variable such that the parameter value is 0.5.
- **`window`** (`int`) –

**`_compute(*idx, **kwargs)`**

Returns: The value of the parameter in the provided index.

**`update(*idx, **kwargs)`**

Updates the value of the parameter in the provided index.

#### Parameters

- **`*idx`** (`list`) – index of the parameter whose number of visits has to be updated.
- **`target`** (`float`) – Value of the target variable;
- **`factor`** (`float`) – Multiplicative factor for the parameter value, useful when the parameter depend on another parameter value.

**`__call__(*idx, **kwargs)`**

Update and return the parameter in the provided index.

**Parameters** `*idx` (`list`) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

**`get_value(*idx, **kwargs)`**

Return the current value of the parameter in the provided index.

**Parameters** `*idx` (`list`) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

#### shape

The shape of the table of parameters.

**Type** Returns

```
class mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter(value,
                                                                                 ex-
                                                                                 po-
                                                                                 nen-
                                                                                 tial=False,
                                                                                 min_value=None,
                                                                                 tol=1.0,
                                                                                 win-
                                                                                 dow=100,
                                                                                 size=(1,
                                                                                 ))
```

Bases: `mushroom_rl.utils.variance_parameters.WindowedVarianceParameter`

Class implementing a parameter that decreases with the target variance, where the variance is computed in a fixed length window.

**`__call__(idx, **kwargs)`**

Update and return the parameter in the provided index.

**Parameters** `*idx`(*list*) – index of the parameter to return.

**Returns** The updated parameter in the provided index.

**`__init__(value, exponential=False, min_value=None, tol=1.0, window=100, size=(1,))`**

Constructor.

**Parameters**

- `tol`(*float*) – value of the variance of the target variable such that the parameter value is 0.5.

- `window`(*int*) –

**`_compute(idx, **kwargs)`**

Returns: The value of the parameter in the provided index.

**`get_value(idx, **kwargs)`**

Return the current value of the parameter in the provided index.

**Parameters** `*idx`(*list*) – index of the parameter to return.

**Returns** The current value of the parameter in the provided index.

**`shape`**

The shape of the table of parameters.

**Type** Returns

**`update(idx, **kwargs)`**

Updates the value of the parameter in the provided index.

**Parameters**

- `*idx`(*list*) – index of the parameter whose number of visits has to be updated.
- `target`(*float*) – Value of the target variable;
- `factor`(*float*) – Multiplicative factor for the parameter value, useful when the parameter depend on another parameter value.

### 3.11.16 Viewer

**`class mushroom_rl.utils.viewer.ImageViewer(size, dt)`**

Bases: `object`

Interface to pygame for visualizing plain images.

**`__init__(size, dt)`**

Constructor.

**Parameters**

- `size`(*[list, tuple]*) – size of the displayed image;
- `dt`(*float*) – duration of a control step.

**`display(img)`**

Display given frame.

**Parameters** `img` – image to display.

---

```
class mushroom_rl.utils.viewer.Viewer(env_width, env_height, width=500, height=500, background=(0, 0, 0))
```

Bases: object

Interface to pygame for visualizing mushroom native environments.

```
__init__(env_width, env_height, width=500, height=500, background=(0, 0, 0))
```

Constructor.

#### Parameters

- **env\_width** (*int*) – The x dimension limit of the desired environment;
- **env\_height** (*int*) – The y dimension limit of the desired environment;
- **width** (*int*, 500) – width of the environment window;
- **height** (*int*, 500) – height of the environment window;
- **background** (*tuple*, (0, 0, 0)) – background color of the screen.

**screen**

Property.

**Returns** The screen created by this viewer.

**size**

Property.

**Returns** The size of the screen.

```
line(start, end, color=(255, 255, 255), width=1)
```

Draw a line on the screen.

#### Parameters

- **start** (*np.ndarray*) – starting point of the line;
- **end** (*np.ndarray*) – end point of the line;
- **color** (*tuple* (255, 255, 255)) – color of the line;
- **width** (*int*, 1) – width of the line.

```
square(center, angle, edge, color=(255, 255, 255), width=0)
```

Draw a square on the screen and apply a roto-translation to it.

#### Parameters

- **center** (*np.ndarray*) – the center of the polygon;
- **angle** (*float*) – the rotation to apply to the polygon;
- **edge** (*float*) – length of an edge;
- **color** (*tuple*, (255, 255, 255)) – the color of the polygon;
- **width** (*int*, 0) – the width of the polygon line, 0 to fill the polygon.

```
polygon(center, angle, points, color=(255, 255, 255), width=0)
```

Draw a polygon on the screen and apply a roto-translation to it.

#### Parameters

- **center** (*np.ndarray*) – the center of the polygon;
- **angle** (*float*) – the rotation to apply to the polygon;
- **points** (*list*) – the points of the polygon w.r.t. the center;

- **color** (*tuple, (255, 255, 255)*) – the color of the polygon;
- **width** (*int, 0*) – the width of the polygon line, 0 to fill the polygon.

**circle** (*center, radius, color=(255, 255, 255), width=0*)

Draw a circle on the screen.

#### Parameters

- **center** (*np.ndarray*) – the center of the circle;
- **radius** (*float*) – the radius of the circle;
- **color** (*tuple, (255, 255, 255)*) – the color of the circle;
- **width** (*int, 0*) – the width of the circle line, 0 to fill the circle.

**arrow\_head** (*center, scale, angle, color=(255, 255, 255)*)

Draw an arrow head.

#### Parameters

- **center** (*np.ndarray*) – the position of the arrow head;
- **scale** (*float*) – scale of the arrow, correspond to the length;
- **angle** (*float*) – the angle of rotation of the arrow head;
- **color** (*tuple, (255, 255, 255)*) – the color of the arrow.

**force\_arrow** (*center, direction, force, max\_force, max\_length, color=(255, 255, 255), width=1*)

Draw a force arrow, i.e. an arrow representing a force. The length of the arrow is directly proportional to the force value.

#### Parameters

- **center** (*np.ndarray*) – the point where the force is applied;
- **direction** (*np.ndarray*) – the direction of the force;
- **force** (*float*) – the applied force value;
- **max\_force** (*float*) – the maximum force value;
- **max\_length** (*float*) – the length to use for the maximum force;
- **color** (*tuple, (255, 255, 255)*) – the color of the arrow;
- **width** (*int, 1*) – the width of the force arrow.

**torque\_arrow** (*center, torque, max\_torque, max\_radius, color=(255, 255, 255), width=1*)

Draw a torque arrow, i.e. a circular arrow representing a torque. The radius of the arrow is directly proportional to the torque value.

#### Parameters

- **center** (*np.ndarray*) – the point where the torque is applied;
- **torque** (*float*) – the applied torque value;
- **max\_torque** (*float*) – the maximum torque value;
- **max\_radius** (*float*) – the radius to use for the maximum torque;
- **color** (*tuple, (255, 255, 255)*) – the color of the arrow;
- **width** (*int, 1*) – the width of the torque arrow.

**background\_image**(*img*)

Use the given image as background for the window, rescaling it appropriately.

**Parameters** *img* – the image to be used.

**function**(*x\_s, x\_e, f, n\_points=100, width=1, color=(255, 255, 255)*)

Draw the graph of a function in the image.

**Parameters**

- **x\_s** (*float*) – starting x coordinate;
- **x\_e** (*float*) – final x coordinate;
- **f** (*function*) – the function that maps x coordinates into y coordinates;
- **n\_points** (*int, 100*) – the number of segments used to approximate the function to draw;
- **width** (*int, 1*) – the width of the line drawn;
- **color** (*tuple, (255, 255, 255)*) – the color of the line.

**display**(*s*)

Display current frame and initialize the next frame to the background color.

**Parameters** *s* – time to wait in visualization.

**close**()

Close the viewer, destroy the window.

## 3.12 How to make a simple experiment

The main purpose of MushroomRL is to simplify the scripting of RL experiments. A standard example of a script to run an experiment in MushroomRL, consists of:

- an **initial part** where the setting of the experiment are specified;
- a **middle part** where the experiment is run;
- a **final part** where operations like evaluation, plot and save can be done.

A RL experiment consists of:

- a **MDP**;
- an **agent**;
- a **core**.

A **MDP** is the problem to be solved by the agent. It contains the function to move the agent in the environment according to the provided action. The MDP can be simply created with:

```
import numpy as np
from sklearn.ensemble import ExtraTreesRegressor

from mushroom_rl.algorithms.value import FQI
from mushroom_rl.core import Core
from mushroom_rl.environments import CarOnHill
from mushroom_rl.policy import EpsGreedy
from mushroom_rl.utils.dataset import compute_J
from mushroom_rl.utils.parameters import Parameter
```

(continues on next page)

(continued from previous page)

```
mdp = CarOnHill()
```

A MushroomRL **agent** is the algorithm that is run to learn in the MDP. It consists of a policy approximator and of the methods to improve the policy during the learning. It also contains the features to extract in the case of MDP with continuous state and action spaces. An agent can be defined this way:

```
# Policy
epsilon = Parameter(value=1.)
pi = EpsGreedy(epsilon=epsilon)

# Approximator
approximator_params = dict(input_shape=mdp.info.observation_space.shape,
                            n_actions=mdp.info.action_space.n,
                            n_estimators=50,
                            min_samples_split=5,
                            min_samples_leaf=2)
approximator = ExtraTreesRegressor

# Agent
agent = FQI(mdp.info, pi, approximator, n_iterations=20,
            approximator_params=approximator_params)
```

This piece of code creates the policy followed by the agent (e.g.  $\epsilon$ -greedy) with  $\epsilon = 1$ . Then, the policy approximator is created specifying the parameters to create it and the class (in this case, the `ExtraTreesRegressor` class of scikit-learn is used). Eventually, the agent is created calling the algorithm class and providing the approximator and the policy, together with parameters used by the algorithm.

To run the experiment, the `core` module has to be used. This module requires the agent and the MDP object and contains the function to learn in the MDP and evaluate the learned policy. It can be created with:

```
core = Core(agent, mdp)
```

Once the core has been created, the agent can be trained collecting a dataset and fitting the policy:

```
core.learn(n_episodes=1000, n_episodes_per_fit=1000)
```

In this case, the agent's policy is fitted only once, after that 1000 episodes have been collected. This is a common practice in batch RL algorithms such as FQI where, initially, samples are randomly collected and then the policy is fitted using the whole dataset of collected samples.

Eventually, some operations to evaluate the learned policy can be done. This way the user can, for instance, compute the performance of the agent through the collected rewards during an evaluation run. Fixing  $\epsilon = 0$ , the greedy policy is applied starting from the provided initial states, then the average cumulative discounted reward is returned.

```
pi.set_epsilon(Parameter(0.))
initial_state = np.array([[-.5, 0.]])
dataset = core.evaluate(initial_states=initial_state)

print(compute_J(dataset, gamma=mdp.info.gamma))
```

### 3.13 How to make an advanced experiment

Continuous MDPs are a challenging class of problems to solve in RL. In these problems, a tabular regressor is not enough to approximate the Q-function, since there are an infinite number of states/actions. The solution to solve them

is to use a function approximator (e.g. neural network) fed with the raw values of states and actions. In the case a linear approximator is used, it is convenient to enlarge the input space with the space of non-linear **features** extracted from the raw values. This way, the linear approximator is often able to solve the MDPs, despite its simplicity. Many RL algorithms rely on the use of a linear approximator to solve a MDP, therefore the use of features is very important. This tutorial shows how to solve a continuous MDP in MushroomRL using an algorithm that requires the use of a linear approximator.

Initially, the MDP and the policy are created:

```
import numpy as np

from mushroom_rl.algorithms.value import SARSALambdaContinuous
from mushroom_rl.approximators.parametric import LinearApproximator
from mushroom_rl.core import Core
from mushroom_rl.features import Features
from mushroom_rl.features.tiles import Tiles
from mushroom_rl.policy import EpsGreedy
from mushroom_rl.utils.callbacks import CollectDataset
from mushroom_rl.utils.parameters import Parameter
from mushroom_rl.environments import Gym

# MDP
mdp = Gym(name='MountainCar-v0', horizon=np.inf, gamma=1.)

# Policy
epsilon = Parameter(value=0.)
pi = EpsGreedy(epsilon=epsilon)
```

This is an environment created with the MushroomRL interface to the OpenAI Gym library. Each environment offered by OpenAI Gym can be created this way simply providing the corresponding id in the `name` parameter, except for the Atari that are managed by a separate class. After the creation of the MDP, the tiles features are created:

```
n_tilings = 10
tilings = Tiles.generate(n_tilings, [10, 10],
                        mdp.info.observation_space.low,
                        mdp.info.observation_space.high)
features = Features(tilings=tilings)

approximator_params = dict(input_shape=(features.size,),
                            output_shape=(mdp.info.action_space.n,),
                            n_actions=mdp.info.action_space.n)
```

In this example, we use sparse coding by means of **tiles** features. The `generate` method generates `n_tilings` grids of 10x10 tilings evenly spaced (the way the tilings are created is explained in “*Reinforcement Learning: An Introduction*”, Sutton & Barto, 1998). Eventually, the grid is passed to the `Features` factory method that returns the `features` class.

MushroomRL offers other type of features such a **radial basis functions** and **polynomial** features. The former have also a faster implementation written in Tensorflow that can be used transparently.

Then, the agent is created as usual, but this time passing the feature to it. It is important to notice that the learning rate is divided by the number of tilings for the correctness of the update (see “*Reinforcement Learning: An Introduction*”, Sutton & Barto, 1998 for details). After that, the learning is run as usual:

```
learning_rate = Parameter(.1 / n_tilings)
```

(continues on next page)

(continued from previous page)

```

agent = SARSALambdaContinuous(mdp.info, pi, LinearApproximator,
                               approximator_params=approximator_params,
                               learning_rate=learning_rate,
                               lambda_coeff=.9, features=features)

# Algorithm
collect_dataset = CollectDataset()
callbacks = [collect_dataset]
core = Core(agent, mdp, callbacks_episode=callbacks)

# Train
core.learn(n_episodes=100, n_steps_per_fit=1)

```

To visualize the learned policy the rendering method of OpenAI Gym is used. To activate the rendering in the environments that supports it, it is necessary to set `render=True`.

```
core.evaluate(n_episodes=1, render=True)
```

## 3.14 How to create a regressor

MushroomRL offers a high-level interface to build function regressors. Indeed, it transparently manages regressors for generic functions and Q-function regressors. The user should not care about the low-level implementation of these regressors and should only use the `Regressor` interface. This interface creates a Q-function regressor or a `GenericRegressor` depending on whether the `n_actions` parameter is provided to the constructor or not.

### 3.14.1 Usage of the Regressor interface

**When the action space of RL problems is finite and the adopted approach is value-based**, we want to compute the Q-function of each action. In MushroomRL, this is possible using:

- a Q-function regressor with a different approximator for each action (`ActionRegressor`);
- a single Q-function regressor with a different output for each action (`QRegressor`).

The `QRegressor` is suggested when the number of discrete actions is high, due to memory reasons.

The user can create create a `QRegressor` or an `ActionRegressor`, setting the `output_shape` parameter of the `Regressor` interface. If it is set to `(1,)`, an `ActionRegressor` is created; otherwise if it is set to the number of discrete actions, a `QRegressor` is created.

### 3.14.2 Example

Initially, the MDP, the policy and the features are created:

```

import numpy as np

from mushroom_rl.algorithms.value import SARSALambdaContinuous
from mushroom_rl.approximators.parametric import LinearApproximator
from mushroom_rl.core import Core
from mushroom_rl.environments import *
from mushroom_rl.features import Features

```

(continues on next page)

(continued from previous page)

```

from mushroom_rl.features.tiles import Tiles
from mushroom_rl.policy import EpsGreedy
from mushroom_rl.utils.callbacks import CollectDataset
from mushroom_rl.utils.parameters import Parameter

# MDP
mdp = Gym(name='MountainCar-v0', horizon=np.inf, gamma=1.)

# Policy
epsilon = Parameter(value=0.)
pi = EpsGreedy(epsilon=epsilon)

# Q-function approximator
n_tilings = 10
tilings = Tiles.generate(n_tilings, [10, 10],
                        mdp.info.observation_space.low,
                        mdp.info.observation_space.high)
features = Features(tilings=tilings)

# Agent
learning_rate = Parameter(.1 / n_tilings)

```

The following snippet, sets the output shape of the regressor to the number of actions, creating a QRegressor:

```

approximator_params = dict(input_shape=(features.size,),
                           output_shape=(mdp.info.action_space.n,),
                           n_actions=mdp.info.action_space.n)

```

If you prefer to use an ActionRegressor, simply set the number of actions to (1,):

```

approximator_params = dict(input_shape=(features.size,),
                           output_shape=(1,),
                           n_actions=mdp.info.action_space.n)

```

Then, the rest of the code fits the approximator and runs the evaluation rendering the behaviour of the agent:

```

agent = SARSAContinuous(mdp.info, pi, LinearApproximator,
                        approximator_params=approximator_params,
                        learning_rate=learning_rate,
                        lambda_coeff=.9, features=features)

# Algorithm
collect_dataset = CollectDataset()
callbacks = [collect_dataset]
core = Core(agent, mdp, callbacks_episode=callbacks)

# Train
core.learn(n_episodes=100, n_steps_per_fit=1)

# Evaluate
core.evaluate(n_episodes=1, render=True)

```

### 3.14.3 Generic regressor

Whenever the `n_actions` parameter is not provided, the `Regressor` interface creates a `GenericRegressor`. This regressor can be used for general purposes and it is more flexible to be used. It is commonly used in policy search algorithms.

#### Example

Create a dataset of points distributed on a line with random gaussian noise.

```
import numpy as np
from matplotlib import pyplot as plt

from mushroom_rl.approximators import Regressor
from mushroom_rl.approximators.parametric import LinearApproximator

x = np.arange(10).reshape(-1, 1)

intercept = 10
noise = np.random.randn(10, 1) * 1
y = 2 * x + intercept + noise
```

To fit the intercept, polynomial features of degree 1 are created by hand:

```
phi = np.concatenate((np.ones(10).reshape(-1, 1), x), axis=1)
```

The regressor is then created and fit (note that `n_actions` is not provided):

```
regressor = Regressor(LinearApproximator,
                      input_shape=(2,),
                      output_shape=(1,))

regressor.fit(phi, y)
```

Eventually, the approximated function of the regressor is plotted together with the target points. Moreover, the weights and the gradient in point 5 of the linear approximator are printed.

```
print('Weights: ' + str(regressor.get_weights()))
print('Gradient: ' + str(regressor.diff(np.array([[5.]])))))

plt.scatter(x, y)
plt.plot(x, regressor.predict(phi))
plt.show()
```

## 3.15 How to make a deep RL experiment

The usual script to run a deep RL experiment does not significantly differ from the one for a shallow RL experiment. This tutorial shows how to solve `Atari` games in MushroomRL using `DQN`, and how to solve `MuJoCo` tasks using `DDPG`. This tutorial will not explain some technicalities that are already described in the previous tutorials, and will only briefly explain how to run deep RL experiments. Be sure to read the previous tutorials before starting this one.

### 3.15.1 Solving Atari with DQN

This script runs the experiment to solve the Atari Breakout game as described in the DQN paper “*Human-level control through deep reinforcement learning*”, Mnih V. et al., 2015). We start creating the neural network to learn the action-value function:

```
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F

from mushroom_rl.algorithms.value import DQN
from mushroom_rl.approximators.parametric import TorchApproximator
from mushroom_rl.core import Core
from mushroom_rl.environments import Atari
from mushroom_rl.policy import EpsGreedy
from mushroom_rl.utils.dataset import compute_metrics
from mushroom_rl.utils.parameters import LinearParameter, Parameter

class Network(nn.Module):
    n_features = 512

    def __init__(self, input_shape, output_shape, **kwargs):
        super().__init__()

        n_input = input_shape[0]
        n_output = output_shape[0]

        self._h1 = nn.Conv2d(n_input, 32, kernel_size=8, stride=4)
        self._h2 = nn.Conv2d(32, 64, kernel_size=4, stride=2)
        self._h3 = nn.Conv2d(64, 64, kernel_size=3, stride=1)
        self._h4 = nn.Linear(3136, self.n_features)
        self._h5 = nn.Linear(self.n_features, n_output)

        nn.init.xavier_uniform_(self._h1.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h2.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h3.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h4.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h5.weight,
                               gain=nn.init.calculate_gain('linear'))

    def forward(self, state, action=None):
        h = F.relu(self._h1(state.float() / 255.))
        h = F.relu(self._h2(h))
        h = F.relu(self._h3(h))
        h = F.relu(self._h4(h.view(-1, 3136)))
        q = self._h5(h)

        if action is None:
            return q
        else:
            q_acted = torch.squeeze(q.gather(1, action.long()))
```

(continues on next page)

(continued from previous page)

```
    return q_acted
```

Note that the forward function may return all the action-values of state, or only the one for the provided action. This network will be used later in the script. Now, we define useful functions, set some hyperparameters, and create the mdp and the policy pi:

```
def print_epoch(epoch):
    print('#####')
    print('Epoch: ', epoch)
    print('-----')

def get_stats(dataset):
    score = compute_metrics(dataset)
    print(('min_reward: %f, max_reward: %f, mean_reward: %f,' +
          ' games_completed: %d' % score))

    return score

scores = list()

optimizer = dict()
optimizer['class'] = optim.Adam
optimizer['params'] = dict(lr=.00025)

# Settings
width = 84
height = 84
history_length = 4
train_frequency = 4
evaluation_frequency = 250000
target_update_frequency = 10000
initial_replay_size = 50000
max_replay_size = 500000
test_samples = 125000
max_steps = 50000000

# MDP
mdp = Atari('BreakoutDeterministic-v4', width, height, ends_at_life=True,
            history_length=history_length, max_no_op_actions=30)

# Policy
epsilon = LinearParameter(value=1.,
                           threshold_value=.1,
                           n=1000000)
epsilon_test = Parameter(value=.05)
epsilon_random = Parameter(value=1)
pi = EpsGreedy(epsilon=epsilon_random)
```

Differently from the literature, we use Adam as the optimizer.

Then, the approximator:

```
# Approximator
input_shape = (history_length, height, width)
```

(continues on next page)

(continued from previous page)

```
approximator_params = dict(
    network=Network,
    input_shape=input_shape,
    output_shape=(mdp.info.action_space.n,),
    n_actions=mdp.info.action_space.n,
    n_features=Network.n_features,
    optimizer=optimizer,
    loss=F.smooth_l1_loss
)

approximator = TorchApproximator
```

Finally, the agent and the core:

```
# Agent
algorithm_params = dict(
    batch_size=32,
    target_update_frequency=target_update_frequency // train_frequency,
    replay_memory=None,
    initial_replay_size=initial_replay_size,
    max_replay_size=max_replay_size
)

agent = DQN(mdp.info, pi, approximator,
            approximator_params=approximator_params,
            **algorithm_params)

# Algorithm
core = Core(agent, mdp)
```

Eventually, the learning loop is performed. As done in literature, learning and evaluation steps are alternated:

```
# RUN

# Fill replay memory with random dataset
print_epoch(0)
core.learn(n_steps=initial_replay_size,
           n_steps_per_fit=initial_replay_size)

# Evaluate initial policy
pi.set_epsilon(epsilon_test)
mdp.set_episode_end(False)
dataset = core.evaluate(n_steps=test_samples)
scores.append(get_stats(dataset))

for n_epoch in range(1, max_steps // evaluation_frequency + 1):
    print_epoch(n_epoch)
    print('- Learning:')
    # learning step
    pi.set_epsilon(epsilon)
    mdp.set_episode_end(True)
    core.learn(n_steps=evaluation_frequency,
               n_steps_per_fit=train_frequency)

    print('- Evaluation:')
    # evaluation step
    pi.set_epsilon(epsilon_test)
```

(continues on next page)

(continued from previous page)

```
mdp.set_episode_end(False)
dataset = core.evaluate(n_steps=test_samples)
scores.append(get_stats(dataset))
```

### 3.15.2 Solving MuJoCo with DDPG

This script runs the experiment to solve the Walker-Stand MuJoCo task, as implemented in [MuJoCo](#). As with DQN, we start creating the neural networks. For DDPG, we need an actor and a critic network:

```
import numpy as np

import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F

from mushroom_rl.algorithms.actor_critic import DDPG
from mushroom_rl.core import Core
from mushroom_rl.environments.dmc_control_env import DMControl
from mushroom_rl.policy import OrnsteinUhlenbeckPolicy
from mushroom_rl.utils.dataset import compute_J


class CriticNetwork(nn.Module):
    def __init__(self, input_shape, output_shape, n_features, **kwargs):
        super().__init__()

        n_input = input_shape[-1]
        n_output = output_shape[0]

        self._h1 = nn.Linear(n_input, n_features)
        self._h2 = nn.Linear(n_features, n_features)
        self._h3 = nn.Linear(n_features, n_output)

        nn.init.xavier_uniform_(self._h1.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h2.weight,
                               gain=nn.init.calculate_gain('relu'))
        nn.init.xavier_uniform_(self._h3.weight,
                               gain=nn.init.calculate_gain('linear'))

    def forward(self, state, action):
        state_action = torch.cat((state.float(), action.float()), dim=1)
        features1 = F.relu(self._h1(state_action))
        features2 = F.relu(self._h2(features1))
        q = self._h3(features2)

        return torch.squeeze(q)

class ActorNetwork(nn.Module):
    def __init__(self, input_shape, output_shape, n_features, **kwargs):
        super(ActorNetwork, self).__init__()

        n_input = input_shape[-1]
```

(continues on next page)

(continued from previous page)

```

n_output = output_shape[0]

self._h1 = nn.Linear(n_input, n_features)
self._h2 = nn.Linear(n_features, n_features)
self._h3 = nn.Linear(n_features, n_output)

nn.init.xavier_uniform_(self._h1.weight,
                       gain=nn.init.calculate_gain('relu'))
nn.init.xavier_uniform_(self._h2.weight,
                       gain=nn.init.calculate_gain('relu'))
nn.init.xavier_uniform_(self._h3.weight,
                       gain=nn.init.calculate_gain('linear'))

def forward(self, state):
    features1 = F.relu(self._h1(torch.squeeze(state, 1).float()))
    features2 = F.relu(self._h2(features1))
    a = self._h3(features2)

    return a

```

We create the mdp, the policy, and set some hyperparameters:

```

# MDP
horizon = 500
gamma = 0.99
gamma_eval = 1.
mdp = DMControl('walker', 'stand', horizon, gamma)

# Policy
policy_class = OrnsteinUhlenbeckPolicy
policy_params = dict(sigma=np.ones(1) * .2, theta=.15, dt=1e-2)

# Settings
initial_replay_size = 500
max_replay_size = 5000
batch_size = 200
n_features = 80
tau = .001

```

Note that the policy is not instantiated in the script, since in DDPG the instantiation is done inside the algorithm constructor.

We create the actor and the critic approximators:

```

# Approximator
actor_input_shape = mdp.info.observation_space.shape
actor_params = dict(network=ActorNetwork,
                     n_features=n_features,
                     input_shape=actor_input_shape,
                     output_shape=mdp.info.action_space.shape)

actor_optimizer = {'class': optim.Adam,
                  'params': {'lr': 1e-5}}

critic_input_shape = (actor_input_shape[0] + mdp.info.action_space.shape[0],)
critic_params = dict(network=CriticNetwork,
                     optimizer={'class': optim.Adam,

```

(continues on next page)

(continued from previous page)

```
'params': {'lr': 1e-3},
loss=F.mse_loss,
n_features=n_features,
input_shape=critic_input_shape,
output_shape=(1,))
```

Finally, we create the agent and the core:

```
# Agent
agent = DDPG(mdp.info, policy_class, policy_params,
              actor_params, actor_optimizer, critic_params,
              batch_size, initial_replay_size, max_replay_size,
              tau)

# Algorithm
core = Core(agent, mdp)
```

As in DQN, we alternate learning and evaluation steps:

```
# Fill the replay memory with random samples
core.learn(n_steps=initial_replay_size, n_steps_per_fit=initial_replay_size)

# RUN
n_epochs = 40
n_steps = 1000
n_steps_test = 2000

dataset = core.evaluate(n_steps=n_steps_test, render=False)
J = compute_J(dataset, gamma_eval)
print('Epoch: 0')
print('J: ', np.mean(J))

for n in range(n_epochs):
    print('Epoch: ', n+1)
    core.learn(n_steps=n_steps, n_steps_per_fit=1)
    dataset = core.evaluate(n_steps=n_steps_test, render=False)
```

---

## Python Module Index

---

### M

mushroom\_rl.environments.inverted\_pendulum,  
mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic,  
    10   mushroom\_rl.environments.lqr, 78  
mushroom\_rl.algorithms.actor\_critic.deep\_mushroom\_rl.environments.mujoco, 79  
    13   mushroom\_rl.environments.puddle\_world,  
mushroom\_rl.algorithms.agent, 6                               82  
mushroom\_rl.algorithms.policy\_search.black\_box\_optimization, 84  
    29   mushroom\_rl.environments.ship\_steering,  
mushroom\_rl.algorithms.policy\_search.policy\_gradient,  
    24   84  
   mushroom\_rl.features.\_implementations.features\_implementations  
mushroom\_rl.algorithms.value.batch\_td,  
    46   90  
mushroom\_rl.algorithms.value.dqn, 50                               mushroom\_rl.features.basis.fourier, 90  
mushroom\_rl.algorithms.value.td, 33                               mushroom\_rl.features.basis.gaussian\_rbf, 91  
mushroom\_rl.approximators.parametric.linear, 58                               mushroom\_rl.features.basis.polynomial, 92  
    57   mushroom\_rl.features.torch\_approximator, 89  
   mushroom\_rl.features.tensors.gaussian\_tensor, 92  
mushroom\_rl.approximators.regressor, 55                               mushroom\_rl.features.tiles.tiles, 93  
mushroom\_rl.core.core, 8                                       mushroom\_rl.policy.deterministic\_policy, 95  
mushroom\_rl.distributions.distribution,  
    60   mushroom\_rl.policy.gaussian\_policy, 96  
mushroom\_rl.distributions.gaussian, 61                               mushroom\_rl.policy.noise\_policy, 103  
mushroom\_rl.environments.atari, 65                               mushroom\_rl.policy.policy, 94  
mushroom\_rl.environments.car\_on\_hill,  
    68   mushroom\_rl.policy.td\_policy, 104  
mushroom\_rl.environments.cart\_pole, 76                               mushroom\_rl.policy.torch\_policy, 107  
mushroom\_rl.environments.dm\_control\_env,  
    69   mushroom\_rl.solvers.car\_on\_hill, 111  
   mushroom\_rl.solvers.dynamic\_programming, 111  
mushroom\_rl.environments.environment, 7                               mushroom\_rl.utils.angles, 112  
mushroom\_rl.environments.finite\_mdp, 70                               mushroom\_rl.utils.callbacks, 113  
mushroom\_rl.environments.generators.grid\_world, 86                       mushroom\_rl.utils.dataset, 114  
   mushroom\_rl.utils.eligibility\_trace, 115  
mushroom\_rl.environments.generators.simple\_chain, 87                       mushroom\_rl.utils.features, 117  
   mushroom\_rl.utils.folder, 117  
mushroom\_rl.environments.generators.taxi, 88                               mushroom\_rl.utils.minibatches, 118  
   mushroom\_rl.utils.numerical\_gradient, 118  
mushroom\_rl.environments.grid\_world, 71                               mushroom\_rl.utils.parameters, 119  
mushroom\_rl.environments.gym\_env, 74

mushroom\_rl.utils.replay\_memory, 121  
mushroom\_rl.utils.spaces, 124  
mushroom\_rl.utils.table, 125  
mushroom\_rl.utils.torch, 126  
mushroom\_rl.utils.value\_functions, 127  
mushroom\_rl.utils.variance\_parameters,  
    128  
mushroom\_rl.utils.viewer, 132

---

## Index

---

### Symbols

\_\_call\_\_(mushroom\_rl.approximators.regressor.Regressor method), 56  
\_\_call\_\_(mushroom\_rl.distributions.distribution.Distribution method), 60  
\_\_call\_\_(mushroom\_rl.distributions.gaussian.GaussianCholeskyDistribution method), 64  
\_\_call\_\_(mushroom\_rl.distributions.gaussian.GaussianDiagonalDistribution method), 63  
\_\_call\_\_(mushroom\_rl.distributions.gaussian.GaussianDistribution method), 62  
\_\_call\_\_(mushroom\_rl.features.basis.fourier.FourierBasis method), 90  
\_\_call\_\_(mushroom\_rl.features.basis.gaussian\_rbf.GaussianRBF method), 91  
\_\_call\_\_(mushroom\_rl.features.basis.polynomial.PolynomialBasis method), 92  
\_\_call\_\_(mushroom\_rl.features.tiles.tiles.Tiles method), 93  
\_\_call\_\_(mushroom\_rl.policy.deterministic\_policy.DeterministicPolicy method), 95  
\_\_call\_\_(mushroom\_rl.policy.gaussian\_policy.AbstractGaussianPolicy method), 96  
\_\_call\_\_(mushroom\_rl.policy.gaussian\_policy.DiagonalGaussianPolicy method), 100  
\_\_call\_\_(mushroom\_rl.policy.gaussian\_policy.GaussianPolicy method), 98  
\_\_call\_\_(mushroom\_rl.policy.gaussian\_policy.StateLogStdGaussianPolicy method), 102  
\_\_call\_\_(mushroom\_rl.policy.gaussian\_policy.StateStdGaussianPolicy method), 101  
\_\_call\_\_(mushroom\_rl.policy.noise\_policy.OrnsteinUhlenbeckPolicy method), 103  
\_\_call\_\_(mushroom\_rl.policy.policy.ParametricPolicy method), 95  
\_\_call\_\_(mushroom\_rl.policy.policy.Policy method), 94  
\_\_call\_\_(mushroom\_rl.policy.td\_policy.Boltzmann method), 106  
\_\_call\_\_(mushroom\_rl.policy.td\_policy.EpsGreedy method), 105  
\_\_call\_\_(mushroom\_rl.policy.td\_policy.Mellowmax method), 107  
\_\_call\_\_(mushroom\_rl.policy.td\_policy.TDPolicy method), 104  
\_\_call\_\_(mushroom\_rl.policy.torch\_policy.GaussianTorchPolicy method), 109  
\_\_call\_\_(mushroom\_rl.policy.torch\_policy.TorchPolicy method), 107  
\_\_call\_\_(mushroom\_rl.utils.callbacks.Callback method), 110  
\_\_call\_\_(mushroom\_rl.utils.callbacks.CollectDataset method), 113  
\_\_call\_\_(mushroom\_rl.utils.callbacks.CollectMaxQ method), 113  
\_\_call\_\_(mushroom\_rl.utils.callbacks.CollectParameters method), 114  
\_\_call\_\_(mushroom\_rl.utils.callbacks.CollectQ method), 113  
\_\_call\_\_(mushroom\_rl.utils.parameters.AdaptiveParameter method), 113  
\_\_call\_\_(mushroom\_rl.utils.parameters.ExponentialParameter method), 121  
\_\_call\_\_(mushroom\_rl.utils.parameters.LinearParameter method), 120  
\_\_call\_\_(mushroom\_rl.utils.parameters.Parameter method), 120  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.VarianceDecreasing method), 119  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.VarianceIncreasing method), 130  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.VarianceParameter method), 129  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.VariancedWindowedVariance method), 128  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.WindowedVariance method), 131  
\_\_call\_\_(mushroom\_rl.utils.variance\_parameters.WindowedVariance method), 131  
\_\_init\_\_(mushroom\_rl.distributions.distribution.Distribution attribute), 61

```

__init__(mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy (mushroom_rl.algorithms.value.td.DoubleQLearning
    attribute), 97
__init__(mushroom_rl.policy.policy.ParametricPolicy __init__() (mushroom_rl.algorithms.value.td.ExpectedSARSA
    attribute), 95
__init__(mushroom_rl.policy.policy.Policy attribute), __init__() (mushroom_rl.algorithms.value.td.QLearning
    94
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 10
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 11
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 12
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 13
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 14
__init__() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CORDACorQ (mushroom_rl.algorithms.value.td.RLearning
    method), 15
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN (mushroom_rl.algorithms.value.td.SARSA_Lambda
    method), 16
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN (mushroom_rl.algorithms.value.td.SARSA_Lambda_Continuous
    method), 17
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN (mushroom_rl.algorithms.value.td.SpeedyQLearning
    method), 18
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN (mushroom_rl.algorithms.value.td.TrueOnlineSARSA_Lambda
    method), 19
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN (mushroom_rl.algorithms.value.td.WeightedQLearning
    method), 20
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD3 (mushroom_rl.approximators.parametric.linear.LinearApproximator
    method), 21
__init__() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD3 (mushroom_rl.approximators.parametric.torch_Approximator
    method), 22
__init__(mushroom_rl.algorithms.agent.Agent __init__() (mushroom_rl.approximators.regressor.Regressor
    method), 23
__init__(mushroom_rl.algorithms.policy_search.black_box_optimization.BFGS (mushroom_rl.core.Core
    method), 24
__init__(mushroom_rl.algorithms.policy_search.black_box_optimization.BFGS (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution
    method), 25
__init__(mushroom_rl.algorithms.policy_search.black_box_optimization.BFGS (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution
    method), 26
__init__(mushroom_rl.algorithms.policy_search.black_box_optimization.BFGS (mushroom_rl.environments.atari.Atari
    method), 27
__init__(mushroom_rl.algorithms.policy_search.policy_gradient.BFGS (mushroom_rl.environments.atari.LazyFrames
    method), 28
__init__(mushroom_rl.algorithms.policy_search.policy_gradient.BFGS (mushroom_rl.environments.atari.MaxAndSkip
    method), 29
__init__(mushroom_rl.algorithms.value.batch_td.DoubleFQI (mushroom_rl.environments.car_on_hill.CarOnHill
    method), 30
__init__(mushroom_rl.algorithms.value.batch_td.FQI (mushroom_rl.environments.cart_pole.CartPole
    method), 31
__init__(mushroom_rl.algorithms.value.batch_td.LSPI (mushroom_rl.environments.dmc_control_env.DMControl
    method), 32
__init__(mushroom_rl.algorithms.value.dqn.AveragedDQN (mushroom_rl.environments.environment.Environment
    method), 33
__init__(mushroom_rl.algorithms.value.dqn.CategoricalDQN (mushroom_rl.environments.environment.MDPInfo
    method), 34
__init__(mushroom_rl.algorithms.value.dqn.DQN (mushroom_rl.environments.environment.MDPInfo
    method), 35
__init__(mushroom_rl.algorithms.value.dqn.DoubleDQN (mushroom_rl.environments.finite_mdp.FiniteMDP
    method), 36
__init__(mushroom_rl.algorithms.value.dqn.DoubleDQN (mushroom_rl.environments.grid_world.AbstractGridWorld
    method), 37)

```

`method), 71`  
`__init__() (mushroom_rl.environments.grid_world.GridWorld __init__() (mushroom_rl.utils.callbacks.CollectMaxQ  
method), 72`  
`method), 113`  
`__init__() (mushroom_rl.environments.grid_world.GridWorld __init__() (mushroom_rl.utils.callbacks.CollectParameters  
method), 73`  
`method), 114`  
`__init__() (mushroom_rl.environments.gym_env.Gym __init__() (mushroom_rl.utils.callbacks.CollectQ  
method), 74`  
`method), 113`  
`__init__() (mushroom_rl.environments.inverted_pendulum.InvertedPendulum __init__() (mushroom_rl.utils.eligibility_trace.AccumulatingTrace  
method), 75`  
`method), 116`  
`__init__() (mushroom_rl.environments.lqr.LQR __init__() (mushroom_rl.utils.eligibility_trace.ReplacingTrace  
method), 78`  
`method), 116`  
`__init__() (mushroom_rl.environments.mujoco.MuJoCo __init__() (mushroom_rl.utils.parameters.AdaptiveParameter  
method), 79`  
`method), 121`  
`__init__() (mushroom_rl.environments.puddle_world.PuddleWorld __init__() (mushroom_rl.utils.parameters.ExponentialParameter  
method), 82`  
`method), 120`  
`__init__() (mushroom_rl.environments.segway.Segway __init__() (mushroom_rl.utils.parameters.LinearParameter  
method), 84`  
`method), 119`  
`__init__() (mushroom_rl.environments.ship_steering.ShipSteering __init__() (mushroom_rl.utils.parameters.Parameter  
method), 85`  
`method), 119`  
`__init__() (mushroom_rl.features.basis.fourier.FourierBasis __init__() (mushroom_rl.utils.replay_memory.PrioritizedReplayMemory  
method), 90`  
`method), 123`  
`__init__() (mushroom_rl.features.basis.gaussian_rbf.GaussianRBF __init__() (mushroom_rl.utils.replay_memory.ReplayMemory  
method), 91`  
`method), 121`  
`__init__() (mushroom_rl.features.basis.polynomial.PolynomialBasis __init__() (mushroom_rl.utils.replay_memory.SumTree  
method), 92`  
`method), 122`  
`__init__() (mushroom_rl.features.tensors.gaussian_tensor.PyTorchGaussianRBF __init__() (mushroom_rl.utils.spaces.Box  
method), 92`  
`method), 124`  
`__init__() (mushroom_rl.features.tiles.tiles.Tiles __init__() (mushroom_rl.utils.spaces.Discrete  
method), 93`  
`method), 124`  
`__init__() (mushroom_rl.policy.deterministic_policy.DeterministicPolicy __init__() (mushroom_rl.utils.table.EnsembleTable  
method), 95`  
`method), 125`  
`__init__() (mushroom_rl.policy.gaussian_policy.DiagonalGaussianPolicy __init__() (mushroom_rl.utils.table.Table  
method), 99`  
`method), 125`  
`__init__() (mushroom_rl.policy.gaussian_policy.GaussianPolicy __init__() (mushroom_rl.utils.variance_parameters.VarianceDecreasing  
method), 98`  
`method), 130`  
`__init__() (mushroom_rl.policy.gaussian_policy.StateLogStdGaussianPolicy __init__() (mushroom_rl.utils.variance_parameters.VarianceIncreasing  
method), 102`  
`method), 129`  
`__init__() (mushroom_rl.policy.gaussian_policy.StateStdGaussianPolicy __init__() (mushroom_rl.utils.variance_parameters.VarianceParameter  
method), 100`  
`method), 128`  
`__init__() (mushroom_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy __init__() (mushroom_rl.utils.variance_parameters.WindowedVariance  
method), 103`  
`method), 132`  
`__init__() (mushroom_rl.policy.td_policy.Boltzmann __init__() (mushroom_rl.utils.variance_parameters.WindowedVariance  
method), 106`  
`method), 131`  
`__init__() (mushroom_rl.policy.td_policy.EpsGreedy __init__() (mushroom_rl.utils.viewer.ImageViewer  
method), 105`  
`method), 132`  
`__init__() (mushroom_rl.policy.td_policy.Mellowmax __init__() (mushroom_rl.utils.viewer.Viewer  
method), 106`  
`method), 133`  
`__init__() (mushroom_rl.policy.td_policy.TDPolicy __add_save_attr() (mush-  
method), 104`  
`room_rl.algorithms.actor_critic.classic_actor_critic.COPDAC_Q  
method), 104`  
`__init__() (mushroom_rl.policy.torch_policy.GaussianTorchPolicy __add_save_attr() (mush-  
method), 109`  
`room_rl.algorithms.actor_critic.classic_actor_critic.COPDAC_Q  
method), 109`  
`__init__() (mushroom_rl.policy.torch_policy.TorchPolicy __add_save_attr() (mush-  
method), 107`  
`room_rl.algorithms.actor_critic.classic_actor_critic.StochasticAC  
method), 107`  
`__init__() (mushroom_rl.utils.callbacks.Callback __add_save_attr() (mush-  
method), 107`

```

room_rl.algorithms.actor_critic.classic_actor_critic.StochasticCriticAlgorithms.value.dqn.AveragedDQN
method), 13                                         (mush- _add_save_attr() (mush-
                                                               _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.A2C   room_rl.algorithms.value.dqn.CategoricalDQN
                                                               method), 15                                         method), 54
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.DDPG   room_rl.algorithms.value.dqn.DQN   method),
                                                               method), 17                                         51
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.DeepAC room_rl.algorithms.value.dqn.DoubleDQN
                                                               method), 14                                         method), 52
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.PPO   room_rl.algorithms.value.td.DoubleQLearning
                                                               method), 24                                         method), 38
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.SAC   room_rl.algorithms.value.td.ExpectedSARSA
                                                               method), 20                                         method), 35
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.TD3   room_rl.algorithms.value.td.QLearning
                                                               method), 19                                         method), 37
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.actor_critic.deep_actor_critic.TRPO  room_rl.algorithms.value.td.RLearning
                                                               method), 22                                         method), 40
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.agent.Agent   room_rl.algorithms.value.td.RQLearning
                                                               method), 6                                         method), 43
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.black_box_optimization.PGDE rl.algorithms.value.td.SARSA   method),
                                                               method), 31                                         33
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.black_box_optimization.REWS rl.algorithms.value.td.SARSALambda
                                                               method), 32                                         method), 34
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.black_box_optimization.RWRL rl.algorithms.value.td.SARSAContinuous
                                                               method), 29                                         method), 44
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.policy_gradient.GPOMDRoom_rl.algorithms.value.td.SpeedyQLearning
                                                               method), 26                                         method), 39
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.policy_gradient.REINFORGErl.algorithms.value.td.TrueOnlineSARSALambda
                                                               method), 25                                         method), 45
                                                               _add_save_attr() (mush- _add_save_attr() (mush-
                                                               room_rl.algorithms.policy_search.policy_gradient.eNAC   room_rl.algorithms.value.td.WeightedQLearning
                                                               method), 28                                         method), 42
                                                               _add_save_attr() (mush- _bound() (mushroom_rl.environments.atari.Atari
                                                               room_rl.algorithms.value.batch_td.DoubleFQI   static method), 68
                                                               method), 48
                                                               _add_save_attr() (mush- _bound() (mushroom_rl.environments.car_on_hill.CarOnHill
                                                               room_rl.algorithms.value.batch_td.FQI   static method), 68
                                                               method), 47
                                                               _add_save_attr() (mush- _bound() (mushroom_rl.environments.cart_pole.CartPole
                                                               room_rl.algorithms.value.batch_td.LSPI   static method), 77
                                                               method), 49
                                                               _add_save_attr() (mush- _bound() (mushroom_rl.environments.dm_control_env.DMControl
                                                               static method), 70
                                                               method), 8
                                                               _add_save_attr() (mush- _bound() (mushroom_rl.environments.environment.Environment
                                                               static method), 8

```

```

_bound() (mushroom_rl.environments.finite_mdp.FiniteMDP      room_rl.algorithms.policy_search.policy_gradient.GPOMDP
    static method), 71                                         method), 26
_bound() (mushroom_rl.environments.grid_world.AbstractGridWorld_end_update()          (mush-
    static method), 72                                         room_rl.algorithms.policy_search.policy_gradient.REINFORCE
_bound() (mushroom_rl.environments.grid_world.GridWorld       _episode_end_update()          (mush-
    static method), 72                                         room_rl.algorithms.policy_search.policy_gradient.eNAC
_bound() (mushroom_rl.environments.grid_world.GridWorldVanHasselt_rl.algorithms.policy_search.policy_gradient.eNAC
    static method), 73                                         method), 25
_bound() (mushroom_rl.environments.gym_env.Gym      _fit() (mushroom_rl.algorithms.value.batch_td.DoubleFQI
    static method), 75                                         method), 48
_bound() (mushroom_rl.environments.inverted_pendulum.InvertedPendulum_rl.algorithms.value.batch_td.FQI
    static method), 76                                         method), 47
_bound() (mushroom_rl.environments.lqr.LQR  static _fit_boosted()          (mush-
    method), 79                                         room_rl.algorithms.value.batch_td.DoubleFQI
_bound() (mushroom_rl.environments.mujoco.MuJoCo        _fit_boosted()          (mush-
    static method), 82                                         room_rl.algorithms.value.batch_td.DoubleFQI
_bound() (mushroom_rl.environments.puddle_world.PuddleWorld _fit_boosted()          (mush-
    static method), 83                                         room_rl.algorithms.value.batch_td.FQI
_bound() (mushroom_rl.environments.segway.Segway      _init_update()          (mush-
    static method), 84                                         room_rl.algorithms.policy_search.policy_gradient.GPOMDP
_bound() (mushroom_rl.environments.ship_steering.ShipSteering _init_update()          (mush-
    static method), 85                                         room_rl.algorithms.policy_search.policy_gradient.REINFORCE
_bound() (mushroom_rl.utils.parameters.ExponentialParameter _init_update()          (mush-
    method), 120                                         room_rl.algorithms.policy_search.policy_gradient.eNAC
_bound() (mushroom_rl.utils.parameters.LinearParameter _init_update()          (mush-
    method), 120                                         room_rl.algorithms.policy_search.policy_gradient.eNAC
_bound() (mushroom_rl.utils.parameters.Parameter      _next_q() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DDPG
    method), 119                                         method), 27
_bound() (mushroom_rl.utils.variance_parameters.VarianceDecreasingParameter
    method), 130                                         _next_q() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.SAC
_bound() (mushroom_rl.utils.variance_parameters.VarianceIncreasingParameter
    method), 129                                         _next_q() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD3
_bound() (mushroom_rl.utils.variance_parameters.VarianceParameter
    method), 118                                         _next_q() (mushroom_rl.algorithms.value.dqn.AveragedDQN
_bound() (mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter
    method), 132                                         _next_q() (mushroom_rl.algorithms.value.dqn.CategoricalDQN
_bound() (mushroom_rl.utils.variance_parameters.WindowedVarianceParameter
    method), 131                                         _next_q() (mushroom_rl.algorithms.value.dqn.DQN
_bound() (mushroom_rl.environments.mujoco.MuJoCo      (mush-
    static method), 80                                         _next_q() (mushroom_rl.algorithms.value.dqn.DoubleDQN
    method), 50                                         method), 51
_bound() (mushroom_rl.features.basis.polynomial.PolynomialBasis _next_q() (mushroom_rl.algorithms.value.td.RQLearning
    static method), 92                                         method), 44
_bound() (mushroom_rl.algorithms.policy_search.policy_gradient.GPOMDP_e_actor_parameters() (mush-
    method), 26                                         room_rl.algorithms.actor_critic.deep_actor_critic.A2C
_bound() (mushroom_rl.algorithms.policy_search.policy_gradient.REINFORCE_actor_parameters() (mush-
    method), 25                                         room_rl.algorithms.actor_critic.deep_actor_critic.DDPG
_bound() (mushroom_rl.algorithms.policy_search.policy_gradient.eNAGmize_actor_parameters() (mush-
    method), 28                                         room_rl.algorithms.actor_critic.deep_actor_critic.DeepAC
_episode_end_update() (mush-                                         method), 14

```

```

__optimize_actor_parameters()      (mush-      method), 24
    room_rl.algorithms.actor_critic.deep_actor_critic.SACt_load()          (mush-
        method), 21                                         room_rl.algorithms.actor_critic.deep_actor_critic.SAC
__optimize_actor_parameters()      (mush-      method), 21
    room_rl.algorithms.actor_critic.deep_actor_critic.TD3t_load()          (mush-
        method), 19                                         room_rl.algorithms.actor_critic.deep_actor_critic.TD3
_parse() (mushroom_rl.algorithms.policy_search.policy_gradient.GPOMDP)  _post_load()          (mush-
    method), 27
_parse() (mushroom_rl.algorithms.policy_search.policy_gradient.REINFORCE) (mush-
    method), 25                                         room_rl.algorithms.actor_critic.deep_actor_critic.TRPO
_parse() (mushroom_rl.algorithms.policy_search.policy_gradient.eNAC)       (mushroom_rl.algorithms.agent.Agent
    method), 28                                         method), 6
_parse() (mushroom_rl.algorithms.value.td.DoubleQLearning) _post_load()          (mush-
    static method), 38                                         room_rl.algorithms.policy_search.black_box_optimization.PGPE
_parse() (mushroom_rl.algorithms.value.td.ExpectedSARSA) _post_load()          (mush-
    static method), 36                                         room_rl.algorithms.policy_search.black_box_optimization.REPS
_parse() (mushroom_rl.algorithms.value.td.QLearning) _post_load()          (mush-
    static method), 37                                         room_rl.algorithms.policy_search.black_box_optimization.RWR
_parse() (mushroom_rl.algorithms.value.td.RLearning) _post_load()          (mush-
    static method), 40                                         room_rl.algorithms.policy_search.black_box_optimization.RWR
_parse() (mushroom_rl.algorithms.value.td.RQLearning) _post_load()          (mush-
    static method), 43                                         room_rl.algorithms.policy_search.black_box_optimization.RWR
_parse() (mushroom_rl.algorithms.value.td.SARSA) _post_load()          (mush-
    static method), 33                                         room_rl.algorithms.policy_search.policy_gradient.GPOMDP
_parse() (mushroom_rl.algorithms.value.td.SARSA_Lambda) _post_load()          (mush-
    static method), 34                                         room_rl.algorithms.policy_search.policy_gradient.REINFORCE
_parse() (mushroom_rl.algorithms.value.td.SARSA_Lambda_Continued) _post_load()          (mush-
    static method), 44                                         room_rl.algorithms.policy_search.policy_gradient.eNAC
_parse() (mushroom_rl.algorithms.value.td.SpeedyQLearning) _post_load()          (mush-
    static method), 39                                         room_rl.algorithms.value.batch_td.DoubleFQI
_parse() (mushroom_rl.algorithms.value.td.TrueOnlineSARSA_Lambda) _post_load()          (mush-
    static method), 46                                         room_rl.algorithms.value.batch_td.DoubleFQI
_parse() (mushroom_rl.algorithms.value.td.WeightedQLearning) _post_load()          (mush-
    static method), 42                                         room_rl.algorithms.value.batch_td.FQI
_post_load()          (mush-      room_rl.algorithms.value.batch_td.FQI
    room_rl.algorithms.actor_critic.classic_actor_critic.COPDActQd), 47
    method), 10                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.batch_td.LSPI
    room_rl.algorithms.actor_critic.classic_actor_critic.StochasticAGd), 49
    method), 12                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.dqn.AveragedDQN
    room_rl.algorithms.actor_critic.classic_actor_critic.StochasticAGd_AVG
    method), 13                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.dqn.CategoricalDQN
    room_rl.algorithms.actor_critic.deep_actor_critic.A2C   method), 54
    method), 15                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.dqn.DQN   method),
    room_rl.algorithms.actor_critic.deep_actor_critic.DDPG  51
    method), 17                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.dqn.DoubleDQN
    room_rl.algorithms.actor_critic.deep_actor_critic.DeepAC method), 52
    method), 14                                         _post_load()          (mush-
_post_load()          (mush-      room_rl.algorithms.value.td.DoubleQLearning
    room_rl.algorithms.actor_critic.deep_actor_critic.PPO   method), 38

```

```

__post_load()                               (mush-
    room_rl.algorithms.value.td.ExpectedSARSA
    method), 36
__post_load()                               (mush-
    room_rl.algorithms.value.td.QLearning
    method), 37
__post_load()                               (mush-
    room_rl.algorithms.value.td.RLearning
    method), 40
__post_load()                               (mush-
    room_rl.algorithms.value.td.RQLearning
    method), 43
__post_load()                               (mush-
    room_rl.algorithms.value.td.SARSA
    method), 33
__post_load()                               (mush-
    room_rl.algorithms.value.td.SARSALambda
    method), 34
__post_load()                               (mush-
    room_rl.algorithms.value.td.SARSALambdaContinuous
    method), 44
__post_load()                               (mush-
    room_rl.algorithms.value.td.SpeedyQLearning
    method), 39
__post_load()                               (mush-
    room_rl.algorithms.value.td.TrueOnlineSARSALambda
    method), 46
__post_load()                               (mush-
    room_rl.algorithms.value.td.WeightedQLearning
    method), 42
__preprocess()                             (mush-
    room_rl.core.core.Core
    method), 9
__preprocess_action()                      (mush-
    room_rl.environments.mujoco.MuJoCo
    method), 80
__simulation_post_step()                  (mush-
    room_rl.environments.mujoco.MuJoCo
    method), 81
__simulation_pre_step()                  (mush-
    room_rl.environments.mujoco.MuJoCo
    method), 81
__step() (mushroom_rl.core.core.Core
    method), 9
__step_finalize()                          (mush-
    room_rl.environments.mujoco.MuJoCo
    method), 81
__step_init()                             (mush-
    room_rl.environments.mujoco.MuJoCo
    method), 80
__step_update()                           (mush-
    room_rl.algorithms.policy_search.policy_gradient.GPOMDP
    method), 26
__step_update()                           (mush-
    room_rl.algorithms.policy_search.policy_gradient.REINFORCE
    method), 25
__step_update()                           (mush-
    room_rl.algorithms.policy_search.policy_gradient.eNAC
    method), 28
__update() (mushroom_rl.algorithms.policy_search.black_box_optimizat
    method), 30
__update() (mushroom_rl.algorithms.policy_search.black_box_optimizat
    method), 32
__update() (mushroom_rl.algorithms.policy_search.black_box_optimizat
    method), 29
__update() (mushroom_rl.algorithms.value.td.DoubleQLearning
    method), 38
__update() (mushroom_rl.algorithms.value.td.ExpectedSARSA
    method), 35
__update() (mushroom_rl.algorithms.value.td.QLearning
    method), 36
__update() (mushroom_rl.algorithms.value.td.RLearning
    method), 40
__update() (mushroom_rl.algorithms.value.td.RQLearning
    method), 43
__update() (mushroom_rl.algorithms.value.td.SARSA
    method), 33
__update() (mushroom_rl.algorithms.value.td.SARSALambda
    method), 34
__update() (mushroom_rl.algorithms.value.td.SARSALambdaContinuous
    method), 44
__update() (mushroom_rl.algorithms.value.td.SpeedyQLearning
    method), 39
__update() (mushroom_rl.algorithms.value.td.TrueOnlineSARSALambda
    method), 45
__update() (mushroom_rl.algorithms.value.td.WeightedQLearning
    method), 41
__update_parameters()                     (mush-
    room_rl.algorithms.policy_search.policy_gradient.GPOMDP
    method), 27
__update_parameters()                     (mush-
    room_rl.algorithms.policy_search.policy_gradient.REINFORCE
    method), 25
__update_parameters()                     (mush-
    room_rl.algorithms.policy_search.policy_gradient.eNAC
    method), 28
__update_target()                          (mush-
    room_rl.algorithms.value.dqn.AveragedDQN
    method), 54
__update_target()                          (mush-
    room_rl.algorithms.value.dqn.CategoricalDQN
    method), 55
__update_target()                          (mush-
    room_rl.algorithms.value.dqn.DQN
    method), 50
__GPOMDP_target()                        (mush-
    room_rl.algorithms.value.dqn.DoubleDQN
    method), 52

```

## A

A2C (class in *mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic*), 14  
AbstractGaussianPolicy (class in *mushroom\_rl.policy.gaussian\_policy*), 96  
AbstractGridWorld (class in *mushroom\_rl.environments.grid\_world*), 71  
AccumulatingTrace (class in *mushroom\_rl.utils.eligibility\_trace*), 116  
AdaptiveParameter (class in *mushroom\_rl.utils.parameters*), 121  
add() (*mushroom\_rl.utils.replay\_memory.PrioritizedReplayMemory*.*method*), 123  
add() (*mushroom\_rl.utils.replay\_memory.ReplayMemory*.*method*), 121  
add() (*mushroom\_rl.utils.replay\_memory.SumTree*.*method*), 122  
Agent (class in *mushroom\_rl.algorithms.agent*), 6  
arrays\_as\_dataset() (in module *mushroom\_rl.utils.dataset*), 114  
arrow\_head() (*mushroom\_rl.utils.viewer.Viewer*.*method*), 134  
Atari (class in *mushroom\_rl.environments.atari*), 67  
AveragedDQN (class in *mushroom\_rl.algorithms.value.dqn*), 53

## B

background\_image() (*mushroom\_rl.utils.viewer.Viewer*.*method*), 134  
bfs() (in module *mushroom\_rl.solvers.car\_on\_hill*), 111  
Boltzmann (class in *mushroom\_rl.policy.td\_policy*), 106  
Box (class in *mushroom\_rl.utils.spaces*), 124

## C

Callback (class in *mushroom\_rl.utils.callbacks*), 113  
CarOnHill (class in *mushroom\_rl.environments.car\_on\_hill*), 68  
CartPole (class in *mushroom\_rl.environments.cart\_pole*), 76  
CategoricalDQN (class in *mushroom\_rl.algorithms.value.dqn*), 54  
check\_collision() (*mushroom\_rl.environments.mujoco.MuJoCo*.*method*), 81  
circle() (*mushroom\_rl.utils.viewer.Viewer*.*method*), 134  
clean() (*mushroom\_rl.utils.callbacks.Callback*.*method*), 113  
close() (*mushroom\_rl.environments.atari.MaxAndSkip*.*method*), 66

close() (*mushroom\_rl.utils.viewer.Viewer*.*method*), 135  
CollectDataset (class in *mushroom\_rl.utils.callbacks*), 113  
CollectMaxQ (class in *mushroom\_rl.utils.callbacks*), 113  
CollectParameters (class in *mushroom\_rl.utils.callbacks*), 114  
CollectQ (class in *mushroom\_rl.utils.callbacks*), 113  
compute\_advantage() (in module *mushroom\_rl.utils.value\_functions*), 127  
compute\_advantage\_montecarlo() (in module *mushroom\_rl.utils.value\_functions*), 127  
compute\_gae() (in module *mushroom\_rl.utils.value\_functions*), 127  
compute\_J() (in module *mushroom\_rl.utils.dataset*), 115  
compute\_metrics() (in module *mushroom\_rl.utils.dataset*), 115  
compute\_mu() (in module *mushroom\_rl.environments.generators.grid\_world*), 87  
compute\_mu() (in module *mushroom\_rl.environments.generators.taxis*), 89  
compute\_probabilities() (in module *mushroom\_rl.environments.generators.grid\_world*), 86  
compute\_probabilities() (in module *mushroom\_rl.environments.generators.simple\_chain*), 87  
compute\_probabilities() (in module *mushroom\_rl.environments.generators.taxis*), 88  
compute\_reward() (in module *mushroom\_rl.environments.generators.grid\_world*), 86  
compute\_reward() (in module *mushroom\_rl.environments.generators.simple\_chain*), 87  
compute\_reward() (in module *mushroom\_rl.environments.generators.taxis*), 88  
COPDAC\_Q (class in *mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic*), 10  
copy() (*mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic.COPDAC*.*method*), 10  
copy() (*mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic.Stochastic*.*method*), 12  
copy() (*mushroom\_rl.algorithms.actor\_critic.classic\_actor\_critic.Stochastic*.*method*), 13  
copy() (*mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic.A2C*.*method*), 15  
copy() (*mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic.DDPG*.*method*), 17  
copy() (*mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic.DeepAC*.*method*)

```

        method), 14
copy () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.PPO.mushroom_rl.algorithms.value.td.TrueOnlineSARSALambda
        method), 24
copy () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.SAG.mushroom_rl.algorithms.value.td.WeightedQLearning
        method), 21
copy () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD.mushroom_rl.core.core), 8
        method), 19
copy () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TRPO.DDPG (class in mushroom_rl.algorithms.actor_critic.deep_actor_critic),
        method), 22
copy () (mushroom_rl.algorithms.agent.Agent method), 6
copy () (mushroom_rl.algorithms.policy_search.black_box_optimization.PGPE.class (class in mushroom_rl.algorithms.actor_critic.deep_actor_critic),
        method), 31
copy () (mushroom_rl.algorithms.policy_search.black_box_optimization.REPS
        method), 32
copy () (mushroom_rl.algorithms.policy_search.black_box_optimization.RWBR.deterministic_policy (class in mushroom_rl.algorithms.actor_critic.deep_actor_critic),
        method), 30
copy () (mushroom_rl.algorithms.policy_search.gradient.eNACom.rl.policy.gaussian_policy), 99
        method), 28
copy () (mushroom_rl.algorithms.policy_search.gradient.GPOMDP.method), 58
        method), 27
copy () (mushroom_rl.algorithms.policy_search.gradient.REINFORCE.REINFORCE), 60
        method), 25
copy () (mushroom_rl.algorithms.value.batch_td.DoubleFQI.diff () (mushroom_rl.approximators.parametric.linear.LinearApproximator,
        method), 48
method), 57
copy () (mushroom_rl.algorithms.value.batch_td.FQI.diff () (mushroom_rl.approximators.regressor.Regressor,
        method), 47
method), 61
copy () (mushroom_rl.algorithms.value.batch_td.LSPI.diff () (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution,
        method), 49
method), 65
copy () (mushroom_rl.algorithms.value.dqn.AveragedDQN.diff () (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution,
        method), 53
method), 64
copy () (mushroom_rl.algorithms.value.dqn.CategoricalDQN.diff () (mushroom_rl.policy.deterministic_policy.DeterministicPolicy,
        method), 55
method), 62
copy () (mushroom_rl.algorithms.value.dqn.DoubleDQN.diff () (mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy,
        method), 52
method), 96
copy () (mushroom_rl.algorithms.value.dqn.DQN.diff () (mushroom_rl.policy.gaussian_policy.DiagonalGaussianPolicy),
        method), 51
method), 97
copy () (mushroom_rl.algorithms.value.td.DoubleQLearning.diff () (mushroom_rl.policy.gaussian_policy.GaussianPolicy,
        method), 38
method), 100
copy () (mushroom_rl.algorithms.value.td.ExpectedSARSA.diff () (mushroom_rl.policy.gaussian_policy.StateLogStdGaussianPolicy),
        method), 36
method), 98
copy () (mushroom_rl.algorithms.value.td.QLearning.diff () (mushroom_rl.policy.gaussian_policy.StateStdGaussianPolicy),
        method), 37
method), 102
copy () (mushroom_rl.algorithms.value.td.RLearning.diff () (mushroom_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy),
        method), 40
method), 101
copy () (mushroom_rl.algorithms.value.td.RQLearning.diff () (mushroom_rl.policy.policy.ParametricPolicy),
        method), 43
method), 94
copy () (mushroom_rl.algorithms.value.td.SARSA.diff_log () (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution,
        method), 33
method), 61
copy () (mushroom_rl.algorithms.value.td.SARSALambda.diff_log () (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution),
        method), 35
method), 64
copy () (mushroom_rl.algorithms.value.td.SARSALambdaContinuous.diff_log () (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution),
        method), 44
method), 63
copy () (mushroom_rl.algorithms.value.td.SpeedyQLearning.method), 63

```

```

diff_log() (mushroom_rl.distributions.gaussian.GaussianDistribution_rl.algorithms.actor_critic.classic_actor_critic.StochasticAC
    method), 62                                         method), 13
diff_log() (mushroom_rl.policy.deterministic_policy.DeterministicPolicy)                               (mush-
    method), 96
diff_log() (mushroom_rl.policy.gaussian_policy.AbstractGaussianPolicy), 15
    method), 97                                         draw_action()                               (mush-
diff_log() (mushroom_rl.policy.gaussian_policy.DiagonalGaussianPolicy.algorithms.actor_critic.deep_actor_critic.A2C
    method), 99                                         method), 17
diff_log() (mushroom_rl.policy.gaussian_policy.GaussianPolicy).action()                           (mush-
    method), 98                                         room_rl.algorithms.actor_critic.deep_actor_critic.DeepAC
diff_log() (mushroom_rl.policy.gaussian_policy.StateLogStdGaussianPolicy), 4
    method), 102                                         draw_action()                               (mush-
diff_log() (mushroom_rl.policy.gaussian_policy.StateStdGaussianPolicy.algorithms.actor_critic.deep_actor_critic.PPO
    method), 100                                         method), 24
diff_log() (mushroom_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy) ()                         (mush-
    method), 104                                         room_rl.algorithms.actor_critic.deep_actor_critic.SAC
diff_log() (mushroom_rl.policy.policy.ParametricPolicy) draw_action()                           (mush-
    method), 94                                         method), 21
Discrete (class in mushroom_rl.utils.spaces), 124
display() (mushroom_rl.utils.viewer.ImageViewer draw_action()                               (mush-
    method), 132                                         room_rl.algorithms.actor_critic.deep_actor_critic.TD3
display() (mushroom_rl.utils.viewer.Viewer method), 135                                         method), 19
draw_action() (mushroom_rl.policy.torch_policy.GaussianTorchPolicy draw_action()                   (mush-
    method), 110                                         room_rl.algorithms.policy_search.black_box_optimization.PGPE
distribution() (mush- draw_action()                               (mush-
    room_rl.distributions.distribution), 60                                         method), 31
distribution() (mush- draw_action()                               (mush-
    room_rl.policy.torch_policy.TorchPolicy draw_action()                   (mush-
    method), 108                                         room_rl.algorithms.policy_search.black_box_optimization.REPS
distribution_t() (mush- draw_action()                               (mush-
    room_rl.policy.torch_policy.GaussianTorchPolicy draw_action()                   (mush-
    method), 110                                         room_rl.algorithms.policy_search.black_box_optimization.RWR
distribution_t() (mush- draw_action()                               (mush-
    room_rl.policy.torch_policy.TorchPolicy draw_action()                   (mush-
    method), 108                                         room_rl.algorithms.policy_search.policy_gradient.eNAC
DMControl (class in mushroom_rl.environments.dm_control_env), 69 draw_action()                   (mush-
    room_rl.environments.dm_control_env), 69                                         method), 28
DoubleDQN (class in mushroom_rl.algorithms.value.dqn), 51 draw_action()                               (mush-
    room_rl.algorithms.value.dqn), 51                                         method), 27
DoubleFQI (class in mushroom_rl.algorithms.value.batch_td), 47 draw_action()                               (mush-
    room_rl.algorithms.value.batch_td), 47                                         method), 25
DoubleQLearning (class in mushroom_rl.algorithms.value.td), 37 draw_action()                               (mush-
    room_rl.algorithms.value.td), 37                                         room_rl.algorithms.policy_search.policy_gradient.REINFORCE
DQN (class in mushroom_rl.algorithms.value.dqn), 50 draw_action()                               (mush-
    room_rl.algorithms.value.dqn), 50                                         method), 48
draw_action() (mush- draw_action()                               (mush-
    room_rl.algorithms.actor_critic.classic_actor_critic.COPDAGOrl.algorithms.value.batch_td.FQI
    method), 10                                         method), 47
draw_action() (mush- draw_action()                               (mush-
    room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.batch_td.LSPI
    method), 12                                         method), 49
draw_action() (mush- draw_action()                               (mush-

```

<code>room_rl.algorithms.value.dqn.AveragedDQN method), 53</code>	<code>room_rl.policy.gaussian_policy.GaussianPolicy method), 99</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.dqn.CategoricalDQN method), 55</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.gaussian_policy.StateLogStdGaussianPolicy method), 103</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.dqn.DoubleDQN method), 52</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.gaussian_policy.StateStdGaussianPolicy method), 101</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.dqn.DQN method), 50</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy method), 103</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.DoubleQLearning method), 38</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.policy.ParametricPolicy method), 95</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.ExpectedSARSA method), 36</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.policy.Policy method), 94</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.QLearning method), 37</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.td_policy.Boltzmann method), 106</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.RLearning method), 41</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.td_policy.EpsGreedy method), 105</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.RQLearning method), 43</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.td_policy.Mellowmax method), 107</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.SARSA method), 33</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.td_policy.TDPolicy method), 105</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.SARSALambda method), 35</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.torch_policy.GaussianTorchPolicy method), 110</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.SARSALambdaContinuous method), 45</code>	<code>draw_action ()</code> (mush- <code>room_rl.policy.torch_policy.TorchPolicy method), 108</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.SpeedyQLearning method), 39</code>	<code>draw_action_t ()</code> (mush- <code>room_rl.policy.torch_policy.GaussianTorchPolicy method), 109</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.TrueOnlineSARSALambda method), 46</code>	<code>draw_action_t ()</code> (mush- <code>room_rl.policy.torch_policy.TorchPolicy method), 108</code>
<code>draw_action ()</code> (mush- <code>room_rl.algorithms.value.td.WeightedQLearning method), 42</code>	<b>E</b>
<code>draw_action ()</code> (mush- <code>room_rl.policy.deterministic_policy.DeterministicPolicy method), 96</code>	<code>EligibilityTrace ()</code> (in module <code>mush- room_rl.utils.eligibility_trace</code> ), 115
<code>draw_action ()</code> (mush- <code>room_rl.policy.gaussian_policy.AbstractGaussianPolicy method), 97</code>	<code>eNAC</code> (class in <code>mush- room_rl.algorithms.policy_search.policy_gradient</code> ), 27
<code>draw_action ()</code> (mush- <code>room_rl.policy.gaussian_policy.DiagonalGaussianPolicy method), 100</code>	<code>EnsembleTable</code> (class in <code>mushroom_rl.utils.table</code> ), 125
<code>draw_action ()</code> (mush-	<code>entropy ()</code> ( <code>mushroom_rl.policy.torch_policy.GaussianTorchPolicy method</code> ), 110
	<code>Entropy ()</code> ( <code>mushroom_rl.policy.torch_policy.TorchPolicy method</code> ), 108

entropy_t ()	(mush- room_rl.policy.torch_policy.GaussianTorchPolicy method), 109	room_rl.algorithms.policy_search.policy_gradient.GPOMDP method), 27
entropy_t ()	(mush- room_rl.policy.torch_policy.TorchPolicy method), 108	episode_start () (mush- room_rl.algorithms.policy_search.policy_gradient.REINFORCE method), 25
Environment	(class in mush- room_rl.environments.environment), 7	episode_start () (mush- room_rl.algorithms.value.batch_td.DoubleFQI method), 48
episode_start ()	(mush- room_rl.algorithms.actor_critic.classic_actor_critic.COPDAGOrl.algorithms.value.batch_td.FQI method), 11	episode_start () (mush- room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.batch_td.LSPI method), 47
episode_start ()	(mush- room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.batch_td.LSPI method), 11	episode_start () (mush- room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.batch_td.LSPI method), 49
episode_start ()	(mush- room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.dqn.AveragedDQN method), 13	episode_start () (mush- room_rl.algorithms.actor_critic.classic_actor_critic.StochasticACrl.algorithms.value.dqn.AveragedDQN method), 53
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.A2C method), 15	episode_start () (mush- room_rl.algorithms.value.dqn.CategoricalDQN method), 55
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.DDPG method), 17	episode_start () (mush- room_rl.algorithms.value.dqn.DoubleDQN method), 52
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.DeepAC method), 14	episode_start () (mush- room_rl.algorithms.value.dqn.DQN method), 51
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.PPO method), 24	episode_start () (mush- room_rl.algorithms.value.td.DoubleQLearning method), 38
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.SAC method), 21	episode_start () (mush- room_rl.algorithms.value.td.ExpectedSARSA method), 36
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.TD3 method), 19	episode_start () (mush- room_rl.algorithms.value.td.QLearning method), 37
episode_start ()	(mush- room_rl.algorithms.actor_critic.deep_actor_critic.TRPO method), 23	episode_start () (mush- room_rl.algorithms.value.td.RLearning method), 41
episode_start ()	(mush- room_rl.algorithms.agent.Agent 6	episode_start () (mush- room_rl.algorithms.value.td.RQLearning method), 43
episode_start ()	(mush- room_rl.algorithms.policy_search.black_box_optimization.PGDRE method), 31	episode_start () (mush- room_rl.algorithms.value.td.SARSA method), 33
episode_start ()	(mush- room_rl.algorithms.policy_search.black_box_optimization.REWS method), 32	episode_start () (mush- room_rl.algorithms.value.td.SARSALambda method), 34
episode_start ()	(mush- room_rl.algorithms.policy_search.black_box_optimization.RWRI method), 30	episode_start () (mush- room_rl.algorithms.value.td.SARSAContinuous method), 44
episode_start ()	(mush- room_rl.algorithms.policy_search.policy_gradient.eNAC method), 28	episode_start () (mush- room_rl.algorithms.value.td.SpeedyQLearning method), 40
episode_start ()	(mush-	episode_start () (mush-

```

    room_rl.algorithms.value.td.TrueOnlineSARSLambda() (mushroom_rl.algorithms.policy_search.policy_gradient.REINFORCE
method), 45
episode_start() (mushroom_rl.algorithms.value.td.WeightedQLearning
method), 42
episodes_length() (in module mushroom_rl
    room_rl.utils.dataset), 114
EpsGreedy (class in mushroom_rl.policy.td_policy),
    105
euler_to_quat() (in module mushroom_rl
    room_rl.utils.angles), 112
evaluate() (mushroom_rl.core.core.Core method), 9
ExpectedSARSA (class in mushroom_rl
    room_rl.algorithms.value.td), 35
ExponentialParameter (class in mushroom_rl
    room_rl.utils.parameters), 120

F
Features() (in module mushroom_rl
    room_rl.features.features), 89
FiniteMDP (class in mushroom_rl
    room_rl.environments.finite_mdp), 70
fit() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CPOMDP)
fit() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.DQN)
fit() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.RL)
fit() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.SARSA)
fit() (mushroom_rl.algorithms.actor_critic.classic_actor_critic.SARSLambda)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DDPG)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DQN)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.SpeedyQLearning)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD)
fit() (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TRPOL)
fit() (mushroom_rl.algorithms.agent.Agent method), 6
fit() (mushroom_rl.algorithms.policy_search.black_box_optimizer.EPODE)
fit() (mushroom_rl.algorithms.policy_search.black_box_optimizer.ETD)
fit() (mushroom_rl.algorithms.policy_search.black_box_optimizer.MARWR)
fit() (mushroom_rl.algorithms.policy_search.black_box_optimizer.WR)
fit() (mushroom_rl.algorithms.policy_search.black_box_optimizer.WTD)
fit() (mushroom_rl.algorithms.policy_search.eNA)
fit() (mushroom_rl.algorithms.policy_search.eNA)
fit() (mushroom_rl.algorithms.policy_search.eNA)
fit() (mushroom_rl.algorithms.policy_search.eNA)
force_symlink() (in module mushroom_rl
    room_rl.utils.table)

```

```

        room_rl.utils.folder), 117
FourierBasis (class in mushroom_rl.features.basis.fourier), 90
FQI (class in mushroom_rl.algorithms.value.batch_td), 46
function() (mushroom_rl.utils.viewer.Viewer method), 135

G
GaussianCholeskyDistribution (class in mushroom_rl.distributions.gaussian), 64
GaussianDiagonalDistribution (class in mushroom_rl.distributions.gaussian), 63
GaussianDistribution (class in mushroom_rl.distributions.gaussian), 61
GaussianPolicy (class in mushroom_rl.policy.gaussian_policy), 97
GaussianRBF (class in mushroom_rl.features.basis.gaussian_rbf), 91
GaussianTorchPolicy (class in mushroom_rl.policy.torch_policy), 109
generate() (mushroom_rl.environments.lqr.LQR static method), 78
generate() (mushroom_rl.features.basis.fourier.FourierBasis static method), 90
generate() (mushroom_rl.features.basis.gaussian_rbf.GaussianRBF static method), 91
generate() (mushroom_rl.features.basis.polynomial.PolynomialBasis static method), 121
generate() (mushroom_rl.features.tensors.gaussian_tensor.PyTorchGaussianTensor static method), 120
generate() (mushroom_rl.features.tiles.tiles.Tiles static method), 93
generate_grid_world() (in module mushroom_rl.environments.generators.grid_world), 86
generate_simple_chain() (in module mushroom_rl.environments.generators.simple_chain), 87
generate_taxi() (in module mushroom_rl.environments.generators.taxi), 88
get() (mushroom_rl.utils.callbacks.Callback method), 113
get() (mushroom_rl.utils.replay_memory.PrioritizedReplayMemory method), 123
get() (mushroom_rl.utils.replay_memory.ReplayMemory method), 121
get() (mushroom_rl.utils.replay_memory.SumTree method), 122
get_action_features() (in module mushroom_rl.features.features), 90
get_collision_force() (mushroom_rl.environments.mujoco.MuJoCo method), 81
get_gradient() (in module mushroom_rl.utils.torch), 126
get_parameters() (mushroom_rl.distributions.distribution.Distribution method), 61
get_parameters() (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution method), 65
get_parameters() (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution method), 63
get_parameters() (mushroom_rl.distributions.gaussian.GaussianDistribution method), 62
get_q() (mushroom_rl.policy.td_policy.Boltzmann method), 106
get_q() (mushroom_rl.policy.td_policy.EpsGreedy method), 105
get_q() (mushroom_rl.policy.td_policy.Mellowmax method), 107
get_q() (mushroom_rl.policy.td_policy.TDPolicy method), 104
get_regressor() (mushroom_rl.policy.deterministic_policy.DeterministicPolicy method), 95
get_value() (mushroom_rl.utils.parameters.ExponentialParameter method), 121
get_value() (mushroom_rl.utils.parameters.LinearParameter method), 120
get_value() (mushroom_rl.utils.parameters.Parameter method), 119
get_value() (mushroom_rl.utils.variance_parameters.VarianceDecreasingParameter method), 130
get_value() (mushroom_rl.utils.variance_parameters.VarianceIncreasingParameter method), 129
get_value() (mushroom_rl.utils.variance_parameters.VarianceParameter method), 128
get_value() (mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter method), 132
get_value() (mushroom_rl.utils.variance_parameters.WindowedVarianceParameter method), 131
get_weights() (in module mushroom_rl.utils.torch), 126
get_weights() (mushroom_rl.approximators.parametric.linear.LinearApproximator method), 58

```

get\_weights() (mush- info (mushroom\_rl.environments.cart\_pole.CartPole room\_rl.approximators.parametric.torch\_approximator.TorchApproximator method), 60 info (mushroom\_rl.environments.dm\_control\_env.DMControl attribute), 70

get\_weights() (mush- info (mushroom\_rl.environments.environment.Environment room\_rl.approximators.regressor.Regressor method), 56 attribute), 8

get\_weights() (mush- info (mushroom\_rl.environments.finite\_mdp.FiniteMDP room\_rl.policy.deterministic\_policy.DeterministicPolicy method), 96 attribute), 71

get\_weights() (mush- info (mushroom\_rl.environments.grid\_world.AbstractGridWorld room\_rl.policy.gaussian\_policy.AbstractGaussianPolicy method), 97 attribute), 73

get\_weights() (mush- info (mushroom\_rl.environments.grid\_world.GridWorldVanHasselt room\_rl.policy.gaussian\_policy.DiagonalGaussianPolicy method), 99 attribute), 74

get\_weights() (mush- info (mushroom\_rl.environments.grid\_world.GridWorldVanHasselt room\_rl.policy.gaussian\_policy.GaussianPolicy method), 98 attribute), 75

get\_weights() (mush- info (mushroom\_rl.environments.inverted\_pendulum.InvertedPendulum room\_rl.policy.gaussian\_policy.StateLogStdGaussianPolicy method), 102 attribute), 76

get\_weights() (mush- info (mushroom\_rl.environments.lqr.LQR attribute), 79 room\_rl.policy.gaussian\_policy.StateLogStdGaussianPolicy mushroom\_rl.environments.mujoco.MuJoCo attribute), 82

get\_weights() (mush- info (mushroom\_rl.environments.puddle\_world.PuddleWorld room\_rl.policy.gaussian\_policy.StateStdGaussianPolicy attribute), 83 method), 101 info (mushroom\_rl.environments.segway.Segway attribute), 84

get\_weights() (mush- info (mushroom\_rl.environments.ship\_steering.ShipSteering room\_rl.policy.noise\_policy.OrnsteinUhlenbeckPolicy method), 103 attribute), 85

get\_weights() (mush- initialized (mush- room\_rl.policy.policy.ParametricPolicy room\_rl.utils.replay\_memory.PrioritizedReplayMemory method), 95 attribute), 123

get\_weights() (mush- initialized (mush- room\_rl.policy.torch\_policy.GaussianTorchPolicy room\_rl.utils.replay\_memory.ReplayMemory attribute), 122 method), 110

get\_weights() (mush- input\_shape (mush- room\_rl.policy.torch\_policy.TorchPolicy room\_rl.approximators.regressor.Regressor attribute), 56 method), 109

GPOMDP (class in mushroom\_rl.algorithms.policy\_search.policy\_gradient), 26 InvertedPendulum (class in mushroom\_rl.environments.inverted\_pendulum), 75

GridWorld (class in mushroom\_rl.environments.grid\_world), 72 is\_absorbing() (mush- room\_rl.environments.mujoco.MuJoCo method), 82

GridWorldVanHasselt (class in mushroom\_rl.environments.grid\_world), 73

Gym (class in mushroom\_rl.environments.gym\_env), 74

**H**

high (mushroom\_rl.utils.spaces.Box attribute), 124

**I**

ImageViewer (class in mushroom\_rl.utils.viewer), 132

info (mushroom\_rl.environments.atari.Atari attribute), 68

info (mushroom\_rl.environments.car\_on\_hill.CarOnHill attribute), 69

**L**

LazyFrames (class in mushroom\_rl.environments.atari), 67

learn() (mushroom\_rl.core.core.Core method), 8

line() (mushroom\_rl.utils.viewer.Viewer method), 133

LinearApproximator (class in mushroom\_rl.approximators.parametric.linear), 57

LinearParameter (class in mushroom\_rl.utils.parameters), 119

```

load() (mushroom_rl.algorithms.actor_critic.classic_actor_critic) COMBO (mushroom_rl.algorithms.value.td.RLearning
    class method), 11
load() (mushroom_rl.algorithms.actor_critic.classic_actor_critic) COIN (mushroom_rl.algorithms.value.td.RLearning
    class method), 41
load() (mushroom_rl.algorithms.actor_critic.classic_actor_critic) StochasticAC (mushroom_rl.algorithms.value.td.RQLearning
    class method), 43
load() (mushroom_rl.algorithms.actor_critic.classic_actor_critic) StratifiedCRA (mushroom_rl.algorithms.value.td.SARSA
    class method), 34
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) A2C (mushroom_rl.algorithms.value.td.SARSA
    class method), 35
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) DDPG (mushroom_rl.algorithms.value.td.SARSA
    class method), 45
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) DeepA (mushroom_rl.algorithms.value.td.SpeedyQLearning
    class method), 40
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) PPO (mushroom_rl.algorithms.value.td.TrueOnlineSARSA
    class method), 46
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) SA (mushroom_rl.algorithms.value.td.WeightedQLearning
    class method), 42
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) TD3 (mushroom_rl.distributions.distribution.Distribution
    class method), 60
load() (mushroom_rl.algorithms.actor_critic.deep_actor_critic) TRPO (mushroom_rl.distributions.gaussian.GaussianCholeskyDistr
    class method), 64
load() (mushroom_rl.algorithms.agent.Agent) log_pdf() (mushroom_rl.distributions.gaussian.GaussianDiagonalDistr
    method), 63
load() (mushroom_rl.algorithms.policy_search.black_box_optimization) (mushroom_rl.distributions.gaussian.GaussianDistribution
    class method), 31
load() (mushroom_rl.algorithms.policy_search.black_box_optimization) REPS (mushroom_rl.policy.torch_policy.GaussianTorchPolicy
    class method), 32
load() (mushroom_rl.algorithms.policy_search.black_box_optimization) RWR (mushroom_rl.environment.RWR)
    class method), 30
load() (mushroom_rl.algorithms.policy_search.policy_gradient.eNA) (mushroom_rl.policy.torch_policy.TorchPolicy
    class method), 29
load() (mushroom_rl.algorithms.policy_search.policy_gradient) (mushroom_rl.environment.LQR)
    class method), 27
load() (mushroom_rl.algorithms.policy_search.policy_gradient) GROMDP (mushroom_rl.utils.spaces.Box attribute), 124
    class method), 27
load() (mushroom_rl.algorithms.policy_search.policy_gradient) LQR (class in mushroom_rl.environments.lqr), 78
load() (mushroom_rl.algorithms.policy_search.policy_gradient) REDNFORGE (mushroom_rl.algorithms.value.batch_td),
    class method), 26
load() (mushroom_rl.algorithms.value.batch_td) DoubleFQI M
    class method), 48
load() (mushroom_rl.algorithms.value.batch_td) FQI max_p (mushroom_rl.utils.replay_memory.SumTree
    class method), 47
load() (mushroom_rl.algorithms.value.batch_td) LSPI max_priority (mushroom_rl.utils.replay_memory.PrioritizedReplayMemory
    class method), 49
load() (mushroom_rl.algorithms.value.dqn) AveragedDQN attribute), 123
    class method), 54
load() (mushroom_rl.algorithms.value.dqn) CategoricalDQN MaxAndSkip (class in mushroom_rl.environment)
    class method), 55
load() (mushroom_rl.algorithms.value.dqn) DoubleDQN room_rl.environments.atari), 65
    class method), 52
load() (mushroom_rl.algorithms.value.dqn) DQN MDPInfo (class in mushroom_rl.environment), 7
    class method), 51
load() (mushroom_rl.algorithms.value.td) DoubleQLearning room_rl.utils.minibatches), 118
    class method), 38
load() (mushroom_rl.algorithms.value.td) ExpectedSARSA room_rl.utils.minibatches), 118
    class method), 36
load() (mushroom_rl.algorithms.value.td) QLearning room_rl.utils.folder), 117
    class method), 37

```

```

mle() (mushroom_rl.distributions.distribution.Distribution) mushroom_rl.environments.generators.simple_chain
      method), 60
      (module), 87
mle() (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution) environments.generators.taxi
      method), 64
      (module), 88
mle() (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution) environments.grid_world
      method), 63
      (module), 71
mle() (mushroom_rl.distributions.gaussian.GaussianDistribution) mushroom_rl.environments.gym_env (mod-
      62
      ule), 74
model(mushroom_rl.approximators.regressor.Regressor) mushroom_rl.environments.inverted_pendulum
      attribute), 56
      (module), 75
model (mushroom_rl.utils.table.EnsembleTable attribute), 126
      mushroom_rl.environments.lqr (module), 78
MuJoCo (class in mushroom_rl.environments.mujoco), 79
      mushroom_rl.environments.puddle_world
mushroom_rl.algorithms.actor_critic.classic_act (module), 182c
      (module), 10
      mushroom_rl.environments.segway (module),
mushroom_rl.algorithms.actor_critic.deep_actor_critic (module), 84
      (module), 13
      mushroom_rl.environments.ship_steering
mushroom_rl.algorithms.agent (module), 6
      (module), 84
mushroom_rl.algorithms.policy_search.blanksbox_implementations._implementations.features_impl
      (module), 29
      (module), 90
mushroom_rl.algorithms.policy_search.polynomialgradientfeatures.basis.fourier (mod-
      24
      ule), 90
mushroom_rl.algorithms.value.batch_td (module), 46
      mushroom_rl.features.basis.gaussian_rbf
      (module), 91
mushroom_rl.algorithms.value.dqn (module), 50
      mushroom_rl.features.basis.polynomial
      (module), 92
mushroom_rl.algorithms.value.td (module), 33
      mushroom_rl.features.features (module), 89
      mushroom_rl.features.tensors.gaussian_tensor
mushroom_rl.approximators.parametric.linear (module), 57
      (module), 92
mushroom_rl.approximators.parametric.torch_apprendinator (module), 103
      (module), 58
      mushroom_rl.policy.deterministic_policy
      (module), 95
mushroom_rl.approximators.regressor (module), 55
      mushroom_rl.policy.gaussian_policy (mod-
      ule), 96
mushroom_rl.core.core (module), 8
      mushroom_rl.policy.noise_policy (module),
mushroom_rl.distributions.distribution (module), 60
      103
      mushroom_rl.policy.policy (module), 94
mushroom_rl.distributions.gaussian (mod-
      ule), 61
      mushroom_rl.policy.td_policy (module), 104
mushroom_rl.environments.atari (module), 65
      mushroom_rl.policy.torch_policy (module),
      107
mushroom_rl.environments.car_on_hill (module), 68
      mushroom_rl.solvers.car_on_hill (module),
      111
mushroom_rl.environments.cart_pole (mod-
      ule), 76
      mushroom_rl.solvers.dynamic_programming
      (module), 111
mushroom_rl.environments.dm_control_env (module), 69
      mushroom_rl.utils.angles (module), 112
mushroom_rl.environments.environment (module), 7
      mushroom_rl.utils.callbacks (module), 113
mushroom_rl.environments.finite_mdp (module), 70
      mushroom_rl.utils.dataset (module), 114
      mushroom_rl.utils.eligibility_trace
      (module), 115
mushroom_rl.environments.generators.gridworld (module), 86
      mushroom_rl.utils.features (module), 117
      mushroom_rl.utils.folder (module), 117
      mushroom_rl.utils.minibatches (module),

```

118  
mushroom\_rl.utils.numerical\_gradient  
(*module*), 118  
mushroom\_rl.utils.parameters (*module*), 119  
mushroom\_rl.utils.replay\_memory (*module*),  
121  
mushroom\_rl.utils.spaces (*module*), 124  
mushroom\_rl.utils.table (*module*), 125  
mushroom\_rl.utils.torch (*module*), 126  
mushroom\_rl.utils.value\_functions (*mod-  
ule*), 127  
mushroom\_rl.utils.variance\_parameters  
(*module*), 128  
mushroom\_rl.utils.viewer (*module*), 132

**N**

n\_actions (*mushroom\_rl.utils.eligibility\_trace.AccumulatingTrace*.  
attribute), 117  
n\_actions (*mushroom\_rl.utils.eligibility\_trace.ReplacingTrace*.  
attribute), 116  
n\_actions (*mushroom\_rl.utils.table.Table* attribute),  
125  
normalize\_angle () (in *module* *mush-  
room\_rl.utils.angles*), 112  
normalize\_angle\_positive () (in *module* *mush-  
room\_rl.utils.angles*), 112  
numerical\_diff\_dist () (in *module* *mush-  
room\_rl.utils.numerical\_gradient*), 118  
numerical\_diff\_policy () (in *module* *mush-  
room\_rl.utils.numerical\_gradient*), 118

**O**

ObservationType (class in *mush-  
room\_rl.environments.mujoco*), 79  
OrnsteinUhlenbeckPolicy (class in *mush-  
room\_rl.policy.noise\_policy*), 103  
output\_shape (*mush-  
room\_rl.approximators.regressor.Regressor*  
attribute), 56

**P**

Parameter (class in *mushroom\_rl.utils.parameters*),  
119  
parameters () (*mush-  
room\_rl.policy.torch\_policy.GaussianTorchPolicy*  
method), 110  
parameters () (*mush-  
room\_rl.policy.torch\_policy.TorchPolicy*  
method), 109  
parameters\_size (*mush-  
room\_rl.distributions.distribution.Distribution*  
attribute), 61

parameters\_size (*mush-  
room\_rl.distributions.gaussian.GaussianCholeskyDistribution*  
attribute), 65  
parameters\_size (*mush-  
room\_rl.distributions.gaussian.GaussianDiagonalDistribution*  
attribute), 63  
parameters\_size (*mush-  
room\_rl.distributions.gaussian.GaussianDistribution*  
attribute), 62  
ParametricPolicy (class in *mush-  
room\_rl.policy.policy*), 94  
parse\_dataset () (in *module* *mush-  
room\_rl.utils.dataset*), 114  
parse\_grid () (in *module* *mush-  
room\_rl.environments.generators.grid\_world*),  
86  
PolicyTracegrid () (in *module* *mush-  
room\_rl.environments.generators.taxi*), 88  
PPO (class in *mushroom\_rl.algorithms.policy\_search.black\_box\_optimization*),  
30  
Policy (class in *mushroom\_rl.policy.policy*), 94  
policy\_iteration () (in *module* *mush-  
room\_rl.solvers.dynamic\_programming*),  
111  
polygon () (*mushroom\_rl.utils.viewer.Viewer* method),  
133  
PolynomialBasis (class in *mush-  
room\_rl.features.basis.polynomial*), 92  
PPO (class in *mushroom\_rl.algorithms.actor\_critic.deep\_actor\_critic*),  
23  
predict () (*mushroom\_rl.approximators.parametric.linear.LinearApproximator*  
method), 57  
predict () (*mushroom\_rl.approximators.parametric.torch\_approximator*  
method), 59  
predict () (*mushroom\_rl.approximators.regressor.Regressor*  
method), 56  
predict () (*mushroom\_rl.utils.eligibility\_trace.AccumulatingTrace*  
method), 117  
predict () (*mushroom\_rl.utils.eligibility\_trace.ReplacingTrace*  
method), 116  
predict () (*mushroom\_rl.utils.table.EnsembleTable*  
method), 126  
predict () (*mushroom\_rl.utils.table.Table* method),  
125  
PrioritizedReplayMemory (class in *mush-  
room\_rl.utils.replay\_memory*), 123  
PuddleWorld (class in *mush-  
room\_rl.environments.puddle\_world*), 82  
PyTorchGaussianRBF (class in *mush-  
room\_rl.features.tensors.gaussian\_tensor*),  
92

**Q**

QLearning (class in room\_rl.algorithms.value.td), 36  
 quat\_to\_euler() (in module room\_rl.utils.angles), 112

**R**

read\_data() (mushroom\_rl.environments.mujoco.MuJoCo method), 81

Regressor (class in room\_rl.approximators.regressor), 55

REINFORCE (class in room\_rl.algorithms.policy\_search.policy\_gradient), 24

render() (mushroom\_rl.environments.atari.MaxAndSkip method), 66

ReplacingTrace (class in room\_rl.utils.eligibility\_trace), 115

ReplayMemory (class in room\_rl.utils.replay\_memory), 121

REPS (class in room\_rl.algorithms.policy\_search.black\_box\_optimization), 31

reset() (mushroom\_rl.approximators.regressor.Regressor method), 56

reset() (mushroom\_rl.core.core.Core method), 9

reset() (mushroom\_rl.environments.atari.Atari method), 67

reset() (mushroom\_rl.environments.atari.MaxAndSkip method), 66

reset() (mushroom\_rl.environments.car\_on\_hill.CarOnHill method), 68

reset() (mushroom\_rl.environments.cart\_pole.CartPole method), 77

reset() (mushroom\_rl.environments.dm\_control\_env.DMControl method), 69

reset() (mushroom\_rl.environments.environment.Environment method), 7

reset() (mushroom\_rl.environments.finite\_mdp.FiniteMDP method), 70

reset() (mushroom\_rl.environments.grid\_world.AbstractGridWorld method), 71

reset() (mushroom\_rl.environments.grid\_world.GridWorld method), 73

reset() (mushroom\_rl.environments.grid\_world.GridWorldVanHasselt method), 74

reset() (mushroom\_rl.environments.gym\_env.Gym method), 74

reset() (mushroom\_rl.environments.inverted\_pendulum.InvertedPendulum method), 76

reset() (mushroom\_rl.environments.lqr.LQR method), 78

mush- reset() (mushroom\_rl.environments.mujoco.MuJoCo method), 80  
 mush- reset() (mushroom\_rl.environments.puddle\_world.PuddleWorld method), 83  
 mush- reset() (mushroom\_rl.environments.segway.Segway method), 84  
 mush- reset() (mushroom\_rl.environments.ship\_steering.ShipSteering method), 85  
 mush- reset() (mushroom\_rl.policy.deterministic\_policy.DeterministicPolicy method), 96  
 mush- reset() (mushroom\_rl.policy.gaussian\_policy.AbstractGaussianPolicy method), 97  
 mush- reset() (mushroom\_rl.policy.gaussian\_policy.DiagonalGaussianPolicy method), 100  
 mush- reset() (mushroom\_rl.policy.gaussian\_policy.GaussianPolicy method), 99  
 mush- reset() (mushroom\_rl.policy.gaussian\_policy.StateLogStdGaussianPolicy method), 103  
 mush- reset() (mushroom\_rl.policy.gaussian\_policy.StateStdGaussianPolicy method), 101  
 mush- reset() (mushroom\_rl.policy.noise\_policy.OrnsteinUhlenbeckPolicy method), 104  
 mush- reset() (mushroom\_rl.policy.policy.ParametricPolicy method), 95  
 mush- reset() (mushroom\_rl.policy.policy.Policy method), 94  
 mush- reset() (mushroom\_rl.policy.td\_policy.Boltzmann method), 106  
 mush- reset() (mushroom\_rl.policy.td\_policy.EpsGreedy method), 105  
 mush- reset() (mushroom\_rl.policy.td\_policy.Mellowmax method), 107  
 mush- reset() (mushroom\_rl.policy.td\_policy.TDPolicy method), 105  
 mush- reset() (mushroom\_rl.policy.torch\_policy.GaussianTorchPolicy method), 110  
 mush- reset() (mushroom\_rl.policy.torch\_policy.TorchPolicy method), 109  
 mush- reset() (mushroom\_rl.utils.eligibility\_trace.AccumulatingTrace method), 116  
 mush- reset() (mushroom\_rl.utils.eligibility\_trace.ReplacingTrace method), 115  
 mush- reset() (mushroom\_rl.utils.replay\_memory.ReplayMemory method), 122  
 mush- reset() (mushroom\_rl.utils.table.EnsembleTable method), 126  
 mush- reward() (mushroom\_rl.environments.mujoco.MuJoCo method), 82  
 RLearning (class in room\_rl.algorithms.value.td), 40  
 RQLearning (class in room\_rl.algorithms.value.td), 42  
 RWR (class in room\_rl.algorithms.policy\_search.black\_box\_optimization),

29

S

```
SAC (class in mushroom- method), 47
    room_rl.algorithms.actor_critic.deep_actor_critic.save () (mushroom_rl.algorithms.value.batch_td.LSPI
    19 method), 49
sample () (mushroom_rl.distributions.distribution.Distribution) (mushroom_rl.algorithms.value.dqn.AveragedDQN
    method), 54
method), 60
sample () (mushroom_rl.distributions.gaussian.GaussianCholeskyDistribution) (mushroom_rl.algorithms.value.dqn.CategoricalDQN
    method), 55
method), 64
sample () (mushroom_rl.distributions.gaussian.GaussianDiagonalDistribution) (mushroom_rl.algorithms.value.dqn.DoubleDQN
    method), 52
method), 63
sample () (mushroom_rl.distributions.gaussian.GaussianDistribution) (mushroom_rl.algorithms.value.dqn.DQN
    method), 51
method), 61
SARSA (class in mushroom_rl.algorithms.value.td), 33 save () (mushroom_rl.algorithms.value.td.DoubleQLearning
SARSALambda (class in mushroom- method), 39
    room_rl.algorithms.value.td), 34 save () (mushroom_rl.algorithms.value.td.ExpectedSARSA
SARSAContinuous (class in mushroom- method), 36
    room_rl.algorithms.value.td), 44 save () (mushroom_rl.algorithms.value.td.QLearning
save () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.CQL), 37
    method), 11 save () (mushroom_rl.algorithms.value.td.RLearning
save () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.SAC), 40
    method), 12 save () (mushroom_rl.algorithms.value.td.RQLearning
save () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.SAC_AVG
    method), 13 save () (mushroom_rl.algorithms.value.td.SARSA
method), 13 save () (mushroom_rl.algorithms.value.td.SARSALambda
save () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.A2C), 34 save () (mushroom_rl.algorithms.value.td.SARSALambdaContinuous
    method), 15 save () (mushroom_rl.algorithms.value.td.SARSAContinuous
method), 15 save () (mushroom_rl.algorithms.value.td.SpeedyQLearning
save () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DDPG), 35
    method), 17 save () (mushroom_rl.algorithms.value.td.TrueOnlineSARSALambda
method), 17 save () (mushroom_rl.algorithms.value.td.WeightedQLearning
save () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DeepA2C), 45
    method), 14 save () (mushroom_rl.algorithms.value.td.SpeedyQLearning
method), 14 save () (mushroom_rl.algorithms.value.td.TD3)
method), 19 screen (mushroom_rl.utils.viewer.Viewer attribute),
save () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TRPO)
    method), 23 seed () (mushroom_rl.environments.atari.Atari
method), 23 save () (mushroom_rl.environments.atari.MaxAndSkip
method), 6 seed () (mushroom_rl.environments.atari.MaxAndSkip
save () (mushroom_rl.algorithms.policy_search.black_box_optimization.DPES)
    method), 31 seed () (mushroom_rl.environments.car_on_hill.CarOnHill
method), 31 save () (mushroom_rl.environments.car_on_hill.CarOnHill)
    method), 32 seed () (mushroom_rl.environments.cart_pole.CartPole
method), 32 save () (mushroom_rl.environments.cart_pole.CartPole)
    method), 30 seed () (mushroom_rl.environments.dm_control_env.DMControl
method), 30 save () (mushroom_rl.environments.dm_control_env.DMControl)
    method), 29 seed () (mushroom_rl.environments.environment.Environment
method), 29 save () (mushroom_rl.environments.environment.Environment)
    method), 27 seed () (mushroom_rl.environments.finite_mdp.FiniteMDP
method), 27 save () (mushroom_rl.environments.finite_mdp.FiniteMDP)
    method), 26 seed () (mushroom_rl.environments.grid_world.AbstractGridWorld
method), 26 save () (mushroom_rl.environments.grid_world.AbstractGridWorld)
    method), 72
```

```

seed() (mushroom_rl.environments.grid_world.GridWorld) set_sigma() (mushroom_rl.policy.gaussian_policy.GaussianPolicy
    method), 73
seed() (mushroom_rl.environments.grid_world.GridWorldVanHasseltMethod) set_std() (mushroom_rl.policy.gaussian_policy.DiagonalGaussianPolicy
    method), 74
seed() (mushroom_rl.environments.gym_env.Gym) set_weights() (in module mushroom_rl.utils.torch),
    method), 75
seed() (mushroom_rl.environments.inverted_pendulum.InvertedPendulum) set_weights() (mush-
    method), 76
seed() (mushroom_rl.environments.lqr.LQR method), room_rl.approximators.parametric.linear.LinearApproximator
    79
seed() (mushroom_rl.environments.mujoco.MuJoCo) set_weights() (mush-
    method), 80
seed() (mushroom_rl.environments.puddle_world.PuddleWorld) set_weights() (mush-
    method), 83
seed() (mushroom_rl.environments.segway.Segway) set_weights() (mush-
    method), 84
seed() (mushroom_rl.environments.ship_steering.ShipSteering) set_weights() (mush-
    method), 85
Segway (class in mushroom_rl.environments.segway), 84
select_first_episodes() (in module mushroom_rl.utils.dataset), 114
select_random_samples() (in module mushroom_rl.utils.dataset), 115
set_beta() (mushroom_rl.policy.td_policy.Boltzmann) set_weights() (mush-
    method), 99
method), 106
set_beta() (mushroom_rl.policy.td_policy.Mellowmax) set_weights() (mush-
    method), 98
method), 107
set_episode_end() (mush-
    room_rl.environments.atari.Atari) set_weights() (mush-
    room_rl.policy.gaussian_policy.StateLogStdGaussianPolicy
    method), 68
68
set_epsilon() (mush-
    room_rl.policy.td_policy.EpsGreedy) set_weights() (mush-
    room_rl.policy.gaussian_policy.StateStdGaussianPolicy
    method), 105
105
set_parameters() (mush-
    room_rl.distributions.distribution.Distribution) set_weights() (mush-
    room_rl.policy.noise_policy.OrnsteinUhlenbeckPolicy
    method), 61
61
set_parameters() (mush-
    room_rl.distributions.gaussian.GaussianCholeskyDistribution) set_weights() (mush-
    room_rl.policy.policy.ParametricPolicy
    method), 65
65
set_parameters() (mush-
    room_rl.distributions.gaussian.GaussianDiagonalDistribution) set_weights() (mush-
    room_rl.policy.torch_policy.GaussianTorchPolicy
    method), 63
63
set_parameters() (mush-
    room_rl.distributions.gaussian.GaussianDistribution) set_weights() (mush-
    room_rl.policy.torch_policy.TorchPolicy
    method), 62
62
set_q() (mushroom_rl.policy.td_policy.Boltzmann) setup() (mushroom_rl.environments.mujoco.MuJoCo
    method), 106
method), 82
set_q() (mushroom_rl.policy.td_policy.EpsGreedy) shape(mushroom_rl.environments.environment.MDPInfo
    method), 105
attribute), 7
set_q() (mushroom_rl.policy.td_policy.Mellowmax) shape(mushroom_rl.utils.eligibility_trace.AccumulatingTrace
    method), 107
attribute), 117
set_q() (mushroom_rl.policy.td_policy.TDPolicy) shape(mushroom_rl.utils.eligibility_trace.ReplacingTrace
    method), 104
attribute), 116

```

```

shape (mushroom_rl.utils.parameters.ExponentialParameter.step () (mushroom_rl.environments.environment.Environment
attribute), 121
shape (mushroom_rl.utils.parameters.LinearParameter step () (mushroom_rl.environments.finite_mdp.FiniteMDP
attribute), 120
method), 71
shape (mushroom_rl.utils.parameters.Parameter step () (mushroom_rl.environments.grid_world.AbstractGridWorld
attribute), 119
method), 72
shape (mushroom_rl.utils.spaces.Box attribute), 124
shape (mushroom_rl.utils.spaces.Discrete attribute), step () (mushroom_rl.environments.grid_world.GridWorld
124
method), 73
shape (mushroom_rl.utils.table.Table attribute), 125
shape (mushroom_rl.utils.variance_parameters.VarianceDepressingParameter step () (mushroom_rl.environments.gym_env.Gym
attribute), 130
method), 74
shape (mushroom_rl.utils.variance_parameters.VarianceIncreasingParameter step () (mushroom_rl.environments.inverted_pendulum.InvertedPendulum
attribute), 129
method), 76
shape (mushroom_rl.utils.variance_parameters.VarianceParameter step () (mushroom_rl.environments.lqr.LQR method),
attribute), 129
78
shape (mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter step () (mushroom_rl.environments.mujoco.MuJoCo
attribute), 132
method), 80
shape (mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter step () (mushroom_rl.environments.puddle_world.PuddleWorld
attribute), 131
method), 83
ShipSteering (class in mush- step () (mushroom_rl.environments.segway.Segway
room_rl.environments.ship_steering), 84
method), 84
shortest_angular_distance () (in module step () (mushroom_rl.environments.ship_steering.ShipSteering
mushroom_rl.utils.angles), 112
method), 85
size (mushroom_rl.environments.environment.MDPInfo StochasticAC (class in mush-
attribute), 7
room_rl.algorithms.actor_critic.classic_actor_critic),
size (mushroom_rl.utils.replay_memory.ReplayMemory StochasticAC_AVG (class in mush-
attribute), 122
room_rl.algorithms.actor_critic.classic_actor_critic),
size (mushroom_rl.utils.replay_memory.SumTree stop () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.COPD
attribute), 122
method), 11
size (mushroom_rl.utils.spaces.Discrete attribute), 124
size (mushroom_rl.utils.viewer.Viewer attribute), 133
solve_car_on_hill () (in module mush- stop () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.Stocha
room_rl.solvers.car_on_hill), 111
method), 12
SpeedyQLearning (class in mush- stop () (mushroom_rl.algorithms.actor_critic.classic_actor_critic.Stocha
room_rl.algorithms.value.td), 39
method), 13
square () (mushroom_rl.utils.viewer.Viewer method), stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.A2C
133
method), 16
StateLogStdGaussianPolicy (class in mush- stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DDPG
room_rl.policy.gaussian_policy), 101
method), 17
StateStdGaussianPolicy (class in mush- stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.DeepAC
room_rl.policy.gaussian_policy), 100
method), 14
step () (in module mushroom_rl.solvers.car_on_hill), stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.PPO
111
method), 24
step () (mushroom_rl.environments.atari.Atari stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.SAC
method), 67
method), 21
step () (mushroom_rl.environments.atari.MaxAndSkip stop () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TD3
method), 65
method), 19
step () (mushroom_rl.environments.car_on_hill.CarOnHill.CarOnHill top () (mushroom_rl.algorithms.actor_critic.deep_actor_critic.TRPO
method), 68
method), 23
step () (mushroom_rl.environments.cart_pole.CartPole stop () (mushroom_rl.algorithms.agent.Agent method),
method), 77
6
method), 70
step () (mushroom_rl.environments.dmc_control_env.DMCControl stop () (mushroom_rl.algorithms.policy_search.black_box_optimization.I
method), 70
method), 31

```

stop () (mushroom\_rl.algorithms.policy\_search.black\_box\_optimization.REPS\_rl.environments.environment.Environment method), 32  
stop () (mushroom\_rl.algorithms.policy\_search.black\_box\_optimization.RWRL\_rl.environments.finite\_mdp.FiniteMDP method), 71  
stop () (mushroom\_rl.algorithms.policy\_search.policy\_gradient.eNA\_gym\_env.Environment.Environment.Environment method), 29  
stop () (mushroom\_rl.algorithms.policy\_search.policy\_gradient.GROMDP\_gym\_env.Environment.Environment.Environment method), 73  
stop () (mushroom\_rl.algorithms.policy\_search.policy\_gradient.REINFORCE\_gym\_env.Environment.Environment.Environment method), 26  
stop () (mushroom\_rl.algorithms.value.batch\_td.DoubleFQL\_top () (mushroom\_rl.environments.gym\_env.Gym method), 48  
stop () (mushroom\_rl.algorithms.value.batch\_td.FQI stop () (mushroom\_rl.environments.inverted\_pendulum.InvertedPendulum method), 47  
stop () (mushroom\_rl.algorithms.value.batch\_td.LSPI stop () (mushroom\_rl.environments.lqr.LQR method), 79  
stop () (mushroom\_rl.algorithms.value.dqn.AveragedDQNstop () (mushroom\_rl.environments.mujoco.MuJoCo method), 54  
stop () (mushroom\_rl.algorithms.value.dqn.CategoricalDQNstop () (mushroom\_rl.environments.puddle\_world.PuddleWorld method), 55  
stop () (mushroom\_rl.algorithms.value.dqn.DoubleDQN stop () (mushroom\_rl.environments.segway.Segway method), 52  
stop () (mushroom\_rl.algorithms.value.dqn.DQN stop () (mushroom\_rl.environments.ship\_steering.ShipSteering method), 51  
stop () (mushroom\_rl.algorithms.value.td.DoubleQLearning\_gymTree (class in mushroom\_rl.utils.replay\_memory), 122  
stop () (mushroom\_rl.algorithms.value.td.ExpectedSARSA T  
stop () (mushroom\_rl.algorithms.value.td.QLearning Table (class in mushroom\_rl.utils.table), 125  
method), 37 TD3 (class in mushroom\_rl.environments.mush-  
stop () (mushroom\_rl.algorithms.value.td.RLearning room\_rl.algorithms.actor\_critic.deep\_actor\_critic),  
method), 41 17  
stop () (mushroom\_rl.algorithms.value.td.RQLearning TDPolicy (class in mushroom\_rl.policy.td\_policy), 104  
method), 44 Tiles (class in mushroom\_rl.features.tiles.tiles), 93  
stop () (mushroom\_rl.algorithms.value.td.SARSA to\_float\_tensor () (in module mushroom\_rl.utils.torch), 127  
method), 34  
stop () (mushroom\_rl.algorithms.value.td.SARSALambda TorchApproximator (class in mushroom\_rl.environments.mush-  
method), 35 room\_rl.approximators.parametric.torch\_approximator),  
stop () (mushroom\_rl.algorithms.value.td.SARSALambdaContinuous8  
method), 45 TorchPolicy (class in mushroom\_rl.environments.mush-  
stop () (mushroom\_rl.algorithms.value.td.SpeedyQLearning room\_rl.policy.torch\_policy), 107  
method), 40 torque\_arrow () (mushroom\_rl.utils.viewer.Viewer  
stop () (mushroom\_rl.algorithms.value.td.TrueOnlineSARSALambda method), 134  
method), 46 total\_p (mushroom\_rl.utils.replay\_memory.SumTree  
stop () (mushroom\_rl.algorithms.value.td.WeightedQLearning attribute), 122  
method), 42 TRPO (class in mushroom\_rl.environments.mush-  
stop () (mushroom\_rl.environments.atari.Atari room\_rl.algorithms.actor\_critic.deep\_actor\_critic),  
method), 68 21  
stop () (mushroom\_rl.environments.car\_on\_hill.CarOnHill TrueOnlineSARSALambda (class in mushroom\_rl.environments.mush-  
method), 69 room\_rl.algorithms.value.td), 45  
stop () (mushroom\_rl.environments.cart\_pole.CartPole U  
method), 77  
stop () (mushroom\_rl.environments.dm\_control\_env.DMControl unifrom\_grid () (in module mushroom\_rl.utils.features), 117  
method), 70

```

unwrapped (mushroom_rl.environments.atari.MaxAndSkip
attribute), 67
update () (mushroom_rl.policy.td_policy.Boltzmann
method), 106
update () (mushroom_rl.policy.td_policy.EpsGreedy
method), 105
update () (mushroom_rl.policy.td_policy.Mellowmax
method), 107
update () (mushroom_rl.utils.eligibility_trace.AccumulatingTrace
method), 116
update () (mushroom_rl.utils.eligibility_traceReplacingTwo
method), 116
update () (mushroom_rl.utils.parameters.ExponentialParameter
method), 121
update () (mushroom_rl.utils.parameters.LinearParameter
method), 120
update () (mushroom_rl.utils.parameters.Parameter
method), 119
update () (mushroom_rl.utils.replay_memory.PrioritizedReplayMemory
method), 123
update () (mushroom_rl.utils.replay_memory.SumTree
method), 122
update () (mushroom_rl.utils.variance_parameters.VarianceDecreasingParameter
method), 130
update () (mushroom_rl.utils.variance_parameters.VarianceIncreasingParameter
method), 129
update () (mushroom_rl.utils.variance_parameters.VarianceParameter
method), 128
update () (mushroom_rl.utils.variance_parameters.WindowedVarianceIncreasingParameter
method), 132
update () (mushroom_rl.utils.variance_parameters.WindowedVarianceParameter
method), 131
use_cuda (mushroom_rl.policy.torch_policy.GaussianTorchPolicy
attribute), 110
use_cuda (mushroom_rl.policy.torch_policy.TorchPolicy
attribute), 109
V
value_iteration () (in module mushroom_
room_rl.solvers.dynamic_programming), 111
VarianceDecreasingParameter (class in mushroom_
room_rl.utils.variance_parameters), 129
VarianceIncreasingParameter (class in mushroom_
room_rl.utils.variance_parameters), 129
VarianceParameter (class in mushroom_
room_rl.utils.variance_parameters), 128
Viewer (class in mushroom_rl.utils.viewer), 132
W
WeightedQLearning (class in mushroom_
room_rl.algorithms.value.td), 41
weights_size (mush-
room_rl.approximators.parametric.linear.LinearApproximator)
Z
zero_grad () (in module mushroom_rl.utils.torch), 126

```