

---

# deer Documentation

*Release 0.4*

**deer contributors**

**Sep 12, 2018**



---

## Contents

---

<b>1</b>	<b>What is new</b>	<b>3</b>
1.1	Version 0.4 . . . . .	3
1.2	Version 0.3 . . . . .	3
1.3	Version 0.2 . . . . .	3
1.4	Future extensions: . . . . .	4
1.5	How should I cite DeeR? . . . . .	4
<b>2</b>	<b>User Guide</b>	<b>5</b>
2.1	Installation . . . . .	5
2.2	Tutorial . . . . .	6
2.3	Examples . . . . .	7
2.4	Development . . . . .	13
<b>3</b>	<b>API reference</b>	<b>15</b>
3.1	Agent . . . . .	15
3.2	Controller . . . . .	18
3.3	Environment . . . . .	22
3.4	Learning algorithms . . . . .	23
3.5	Policies . . . . .	28
<b>4</b>	<b>Indices and tables</b>	<b>31</b>
	<b>Python Module Index</b>	<b>33</b>



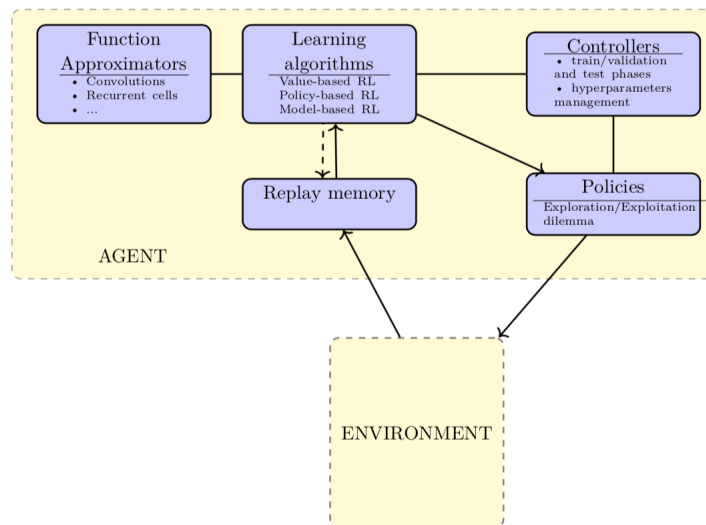
DeeR (Deep Reinforcement) is a python library to train an agent how to behave in a given environment so as to maximize a cumulative sum of rewards (see [What is deep reinforcement learning?](#)).

Here are key advantages of the library:

- You have access within a single library to techniques such as Double Q-learning, prioritized Experience Replay, Deep deterministic policy gradient (DDPG), Combined Reinforcement via Abstract Representations (CRAR), etc.
- This package provides a general framework where observations are made up of any number of elements (scalars, vectors or frames).
- You can easily add up a validation phase that allows to stop the training process before overfitting. This possibility is useful when the environment is dependent on scarce data (e.g. limited time series).

In addition, the framework is made in such a way that it is easy to

- build any environment
- modify any part of the learning process
- use your favorite python-based framework to code your own learning algorithm or neural network architecture. The provided learning algorithms and neural network architectures are based on Keras.



**Figure** General schema of the different elements available in DeeR.

It is a work in progress and input is welcome. Please submit any contribution via pull request.



### 1.1 Version 0.4

- Integration of CRAR that allows to combine the model-free and the model-based approaches via abstract representations.
- Augmented documentation and some interfaces have been updated.

### 1.2 Version 0.3

- Integration of different exploration/exploitation policies and possibility to easily built your own.
- Integration of DDPG for continuous action spaces (see actor-critic)
- *Naming convention for this project* and some interfaces have been updated. This may cause broken backward compatibility. In that case, make the changes to the new convention by looking at the API in this documentation or by looking at the current version of the examples.
- Additional automated tests

### 1.3 Version 0.2

- Standalone python package (you can simply do `pip install deer`)
- Integration of new examples environments : `toy_env_pendulum`, `PLE` and *Gym environment*
- Double Q-learning and prioritized Experience Replay
- Augmented documentation
- First automated tests

## 1.4 Future extensions:

- Several agents interacting in the same environment
- ...

## 1.5 How should I cite DeeR?

Please cite DeeR in your publications if you use it in your research. Here is an example BibTeX entry:

```
@misc{franccoislavet2016deer,  
  title={DeeR},  
  author={Fran\c{c}ois-Lavet, Vincent and others},  
  year={2016},  
  howpublished={\url{https://deer.readthedocs.io/}},  
}
```



## 2.1 Installation

### 2.1.1 Dependencies

This framework is tested to work under Python 3.6.

The required dependencies are NumPy  $\geq 1.10$ , joblib  $\geq 0.9$ . You also need keras or you can write your own learning algorithms using your favorite deep learning framework.

For running some of the examples, Matplotlib  $\geq 1.1.1$  is required. You also sometimes need to install specific dependencies (e.g. for the atari games, you need to install ALE  $\geq 0.4$ ).

We recommend to use the bleeding-edge version and to install it by following the [Developer install instructions](#). If you want a simpler installation procedure and do not intend to modify yourself the learning algorithms etc., you can look at the [User install instructions](#).

### 2.1.2 Developer install instructions

As a developer, you can set you up with the bleeding-edge version of DeeR with:

```
git clone -b master https://github.com/VinF/deer.git
```

Assuming you already have a python environment with `pip`, you can automatically install all the dependencies (except specific dependencies that you may need for some examples) with:

```
pip install -r requirements.txt
```

And you can install the framework as a package using the mode `develop` so that you can make modifications and test without having to re-install the package.

```
python setup.py develop
```

### 2.1.3 User install instructions

You can install the framework with pip:

```
pip install deer
```

For the bleeding edge version (recommended), you can simply use

```
pip install git+git://github.com/VINF/deer.git@master
```

## 2.2 Tutorial

### 2.2.1 What is deep reinforcement learning?

Deep reinforcement learning is the combination of two fields:

- *Reinforcement learning (RL)* is a theory that allows an agent to learn a strategy so as to maximize a sum of cumulated (delayed) rewards from any given environment. If you are not familiar with RL, you can get up to speed easily with by Sutton and Barto.
- *Deep learning* is a branch of machine learning for regression and classification. It is particularly well suited to model high-level abstractions in data by using multiple processing layers composed of multiple non-linear transformations.

This combination allows to learn complex tasks such as playing ATARI games from high-dimensional sensory inputs. For more informations, you can refer to one of the main papers in the domain : .

### 2.2.2 How can I get started?

First, make sure you have installed the package properly by following the steps described in [Installation](#).

The general idea of this framework is that you need to instantiate an agent (along with a learning algorithm) and an environment. In order to perform an experiment, you also need to attach to the agent some controllers for controlling the training and the various parameters of your agent.

The environment should be built specifically for any specific task while learning algorithms (such as q-networks) and many controllers are provided within this package.

The best to get started is to have a look at the [Examples](#) and in particular the two first environments that are simple to understand:

- *Toy environment with time series*
- `toy_env_pendulum`

If you find something that is not yet implemented and if you wish to contribute, you can check the section [Development](#).

### 2.2.3 Any Question?

You can raise questions about the DeeR project on github :

## 2.3 Examples

You can find these examples at the . For each example at least two files are provided:

- A launcher file (whose name usually starts by `run_`).
- An environnement file (whose name usually ends by `_env`).

The launcher file performs different actions:

- It instantiates the environment and the agent along with a learning algorithm (such as a q-network).
- It binds controllers to the agent
- it finally runs the experiment

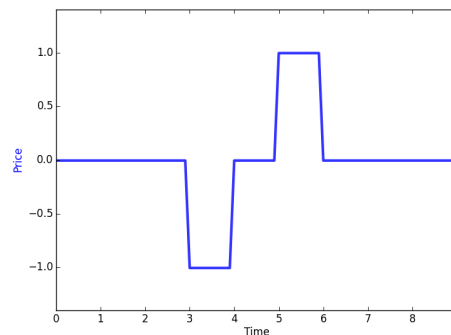
You can get started with the following examples:

### 2.3.1 Toy environment with time series

#### Description of the environnement

This environment simulates the possibility of buying or selling a good. The agent can either have one unit or zero unit of that good. At each transaction with the market, the agent obtains a reward equivalent to the price of the good when selling it and the opposite when buying. In addition, a penalty of 0.5 (negative reward) is added for each transaction.

The price pattern is made by repeating the following signal plus a random constant between 0 and 3:



You can see how this environnement is built by looking into the file `Toy_env.py` in . It is important to note that any environment derives from the base class `Environment` and you can refer to it in order to understand the required methods and their usage.

#### How to run

A minimalist way of running this example can be found in the file `run_toy_env_simple.py` in .

- First, we need to import the agent, the Q-network, the environment and some controllers

```

1 from deer.agent import NeuralAgent
2 from deer.learning_algos.q_net_keras import MyQNetwork
3 from Toy_env import MyEnv as Toy_env
4 import deer.experiment.base_controllers as bc
5
6
```

- Then we instantiate the different elements as follows:

```

1  rng = np.random.RandomState(123456)
2
3  # --- Instantiate environment ---
4  env = Toy_env(rng)
5
6  # --- Instantiate qnetwork ---
7  qnetwork = MyQNetwork(
8      environment=env,
9      random_state=rng)
10
11 # --- Instantiate agent ---
12 agent = NeuralAgent(
13     env,
14     qnetwork,
15     random_state=rng)
16
17 # --- Bind controllers to the agent ---
18 # Before every training epoch, we want to print a summary of the agent's epsilon,
19 ↪ discount and
20 # learning rate as well as the training epoch number.
21 agent.attach(bc.VerboseController())
22
23 # During training epochs, we want to train the agent after every action it takes.
24 # Plus, we also want to display after each training episode (!= than after every
25 ↪ training) the average bellman
26 # residual and the average of the V values obtained during the last episode.
27 agent.attach(bc.TrainerController())
28
29 # All previous controllers control the agent during the epochs it goes through.
30 ↪ However, we want to interleave a
31 # "test epoch" between each training epoch. We do not want these test epoch to
32 ↪ interfere with the training of the
33 # agent. Therefore, we will disable these controllers for the whole duration of the
34 ↪ test epochs interleaved this
35 # way, using the controllersToDisable argument of the InterleavedTestEpochController.
36 ↪ The value of this argument
37 # is a list of the indexes of all controllers to disable, their index reflecting in
38 ↪ which order they were added.
39 agent.attach(bc.InterleavedTestEpochController(
40     epoch_length=500,
41     controllers_to_disable=[0, 1]))
42
43 # --- Run the experiment ---
44 agent.run(n_epochs=100, epoch_length=1000)

```

## Results

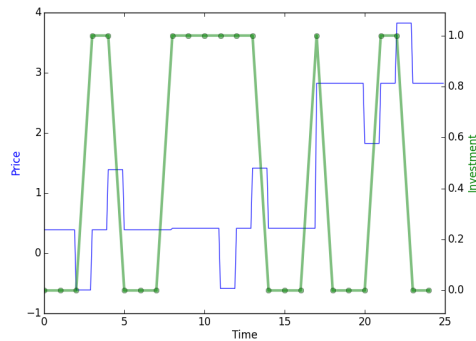
Navigate to the folder `examples/toy_env/` in a terminal window. The example can then be run by using

```
python run_toy_env_simple.py
```

You can also choose the full version of the launcher that specifies the hyperparameters for better performance.

```
python run_toy_env.py
```

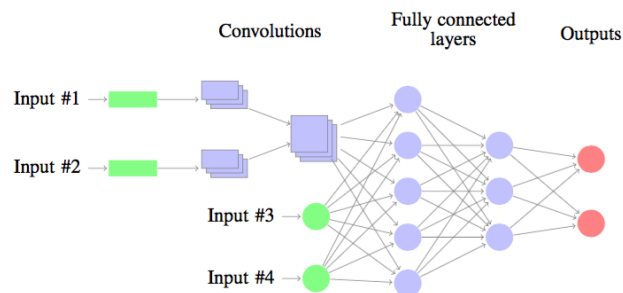
Every 10 epochs, a graph is saved in the ‘toy\_env’ folder. You can then see that kind of behaviour for the test policy at the end of the training:



In this graph, you can see that the agent has successfully learned to take advantage of the price pattern to buy when it is low and to sell when it is high. This example is of course easy due to the fact that the patterns are very systematic which allows the agent to successfully learn it. It is important to note that the results shown are made on a validation set that is different from the training and we can see that learning generalizes well. For instance, the action of buying at time step 7 and 16 is the expected result because in average this will allow to make profit since the agent has no information on the future.

## Using Convolutions VS LSTM's

So far, the neural network was build by using a convolutional architecture as follows:



The neural network processes time series thanks to a set of convolutions layers. The output of the convolutions as well as the other inputs are followed by fully connected layers and the output layer.

When working with deep reinforcement learning, it is also possible to work with LSTM's (see for instance [this](#) )

If you want to use LSTM's architecture, you can import the following libraries

```
from deer.learning_algos.NN_keras_LSTM import NN as NN_keras
```

and then instantiate the qnetwork by specifying the ‘neural\_network’ as follows:

```
qnetwork = MyQNetwork(
    env,
    neural_network=NN_keras)
```

### 2.3.2 Gym environment

Some examples are also provided with the .

Here is the resulting policy for the mountain car example:

Here is the resulting policy for the pendulum example:

### 2.3.3 Two storage devices environment

#### Description of the environment

This example simulates the operation of a realistic micro-grid (such as a smart home for instance) that is not connected to the main utility grid (off-grid) and that is provided with PV panels, batteries and hydrogen storage. The battery has the advantage that it is not limited in instantaneous power that it can provide or store. The hydrogen storage has the advantage that it can store very large quantity of energy.

```
python run_MG_two_storage_devices
```

This example uses the environment defined in `MG_two_storage_devices_env.py`. The agent can either choose to store in the long term storage or take energy out of it while the short term storage handle at best the lack or surplus of energy by discharging itself or charging itself respectively. Whenever the short term storage is empty and cannot handle the net demand a penalty (negative reward) is obtained equal to the value of loss load set to 2euro/kWh.

The state of the agent is made up of an history of two to four punctual observations:

- Charging state of the short term storage (0 is empty, 1 is full)
- Production and consumption (0 is no production or consumption, 1 is maximal production or consumption)
- (Distance to equinox)
- (Predictions of future production : average of the production for the next 24 hours and 48 hours)

Two actions are possible for the agent:

- Action 0 corresponds to discharging the long-term storage
- Action 1 corresponds to charging the long-term storage

**More information can be found in** [Deep Reinforcement Learning Solutions for Energy Microgrids Management](#), Vincent François-Lavet, David Taralla, Damien Ernst, Raphael Fonteneau

#### Annex to the paper

#### PV production and consumption profiles

Solar irradiance varies throughout the year depending on the seasons, and it also varies throughout the day depending on the weather and the position of the sun in the sky relative to the PV panels. The main distinction between these profiles is the difference between summer and winter PV production. In particular, production varies with a factor 1:5 between winter and summer as can be seen from the measurements of PV panels production for a residential customer located in Belgium in the figures below.

A simple residential consumption profile is considered with a daily average consumption of 18kWh (see figure below).

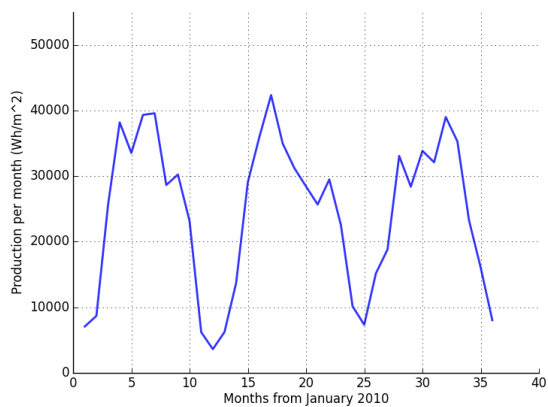


Fig. 1: Total energy produced per month

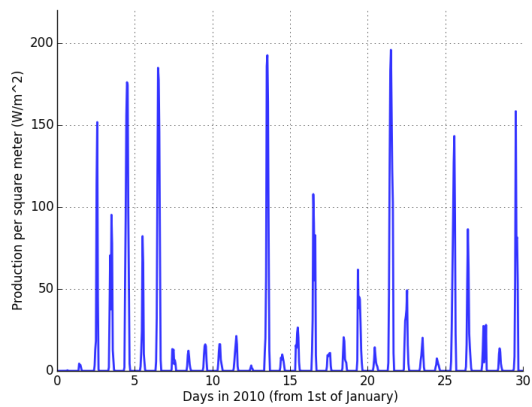


Fig. 2: Typical production in winter

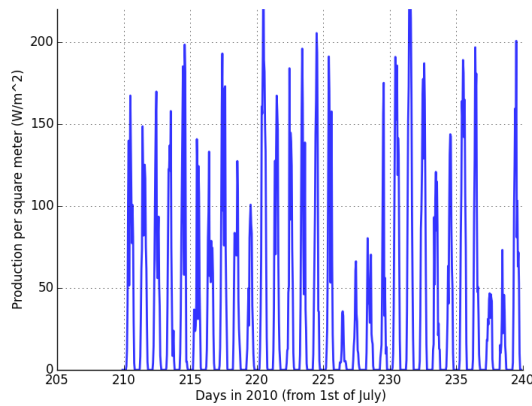


Fig. 3: Typical production in summer

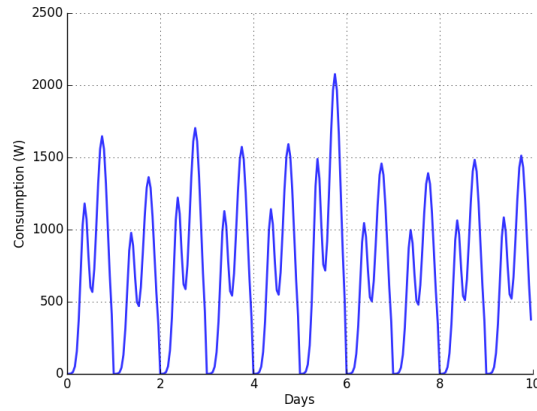


Fig. 4: Representative residential consumption profile

## Main microgrid parameters

Table 1: Data used for the PV panels

cost	$c^{PV}$	$1euro/W_p$
Efficiency	$\eta^{PV}$	18%
Life time	$L^{PV}$	20years

Table 2: Data used for the  $LiFePO_4$  battery

cost	$c^B$	$500euro/kWh$
discharge efficiency	$\eta_0^B$	90%
charge efficiency	$\zeta_0^B$	90%
Maximum instantaneous power	$P^B$	$> 10kW$
Life time	$L^B$	20years

Table 3: Data used for the Hydrogen storage device

cost	$c^{H_2}$	$14euro/W_p$
discharge efficiency	$\eta_0^{H_2}$	65%
charge efficiency	$\zeta_0^{H_2}$	65%
Life time	$L^{H_2}$	20years

Table 4: Data used for reward function

cost endured per kWh not supplied within the microgrid	$k$	$2euro/kWh$
revenue/cost per kWh of hydrogen produced/used	$k^{H_2}$	$0.1euro/kWh$

## 2.3.4 Tasks with planning

You can find the following environments that demonstrate the possibilities of combining model-based and model-free: and .



### 2.3.5 ALE environment

This environment is an interface with the that simulates any ATARI game.

Related paper: Mnih, Volodymyr, et al. “Human-level control through deep reinforcement learning.” Nature 518.7540 (2015): 529-533. (Hyper-parameters tuning is necessary if you want to try to replicate close performances.)

## 2.4 Development

DeeR is a work in progress and contributions are welcome via pull request.

For more information, you can check out this link : .

You should also make sure that you install the repository appropriately for development (see *Developer install instructions*).

### 2.4.1 Guidelines for this project

Here are a few guidelines for this project.

- **Simplicity:** Be easy to use but also easy to understand when one digs into the code. Any additional code should be justified by the usefulness of the feature.
- **Modularity:** The user should be able to easily use its own code with any part of the deer framework (probably at the exception of the core of agent.py that is coded in a very general way).

These guidelines come of course in addition to all good practices for open source development.

### 2.4.2 Naming convention for this project

- All classes and methods have word boundaries using medial capitalization. Classes are written with UpperCamelCase and methods are written with lowerCamelCase respectively. Example: “two words” is rendered as “TwoWords” for the UpperCamelCase (classes) and “twoWords” for the lowerCamelCase (methods).
- All attributes and variables have words separated by underscores. Example: “two words” is rendered as “two\_words”
- If a variable is intended to be ‘private’, it is prefixed by an underscore.



If you are looking for information on a specific function, class or method, this API is for you.

### 3.1 Agent

This module contains classes used to define the standard behavior of the agent. It relies on the controllers, the chosen training/test policy and the learning algorithm to specify its behavior in the environment.

<i>NeuralAgent</i> (environment, learning_algo[, ...])	The NeuralAgent class wraps a learning algorithm (such as a deep Q-network) for training and testing in a given environment.
<i>DataSet</i> (env[, random_state, max_size, ...])	A replay memory consisting of circular buffers for observations, actions, rewards and terminals.

```
class deer.agent.NeuralAgent (environment,    learning_algo,    replay_memory_size=1000000,
                                replay_start_size=None,          batch_size=32,          ran-
                                dom_state=<mtrand.RandomState  object>,    exp_priority=0,
                                train_policy=None, test_policy=None, only_full_history=True)
```

The NeuralAgent class wraps a learning algorithm (such as a deep Q-network) for training and testing in a given environment.

Attach controllers to it in order to conduct an experiment (when to train the agent, when to test,...).

**environment** [object from class Environment] The environment in which the agent interacts

**learning\_algo** [object from class LearningAlgo] The learning algorithm associated to the agent

**replay\_memory\_size** [int] Size of the replay memory. Default : 1000000

**replay\_start\_size** [int] Number of observations (=number of time steps taken) in the replay memory before starting learning. Default: minimum possible according to environment.inputDimensions().

**batch\_size** [int] Number of tuples taken into account for each iteration of gradient descent. Default : 32

**random\_state** [numpy random number generator] Default : random seed.

**exp\_priority** [float] The exponent that determines how much prioritization is used, default is 0 (uniform priority). One may check out Schaul et al. (2016) - Prioritized Experience Replay.

**train\_policy** [object from class Policy] Policy followed when in training mode (mode -1)

**test\_policy** [object from class Policy] Policy followed when in other modes than training (validation and test modes)

**only\_full\_history** [boolean] Whether we wish to train the neural network only on full histories or we wish to fill with zeroes the observations before the beginning of the episode

**avgBellmanResidual** ()  
Returns the average training loss on the epoch

**avgEpisodeVValue** ()  
Returns the average V value on the episode (on time steps where a non-random action has been taken)

**discountFactor** ()  
Get the discount factor

**dumpNetwork** (*fname*, *nEpoch=-1*)  
Dump the network

**fname** [string] Name of the file where the network will be dumped

**nEpoch** [int] Epoch number (Optional)

**learningRate** ()  
Get the learning rate

**overrideNextAction** (*action*)  
Possibility to override the chosen action. This possibility should be used on the signal OnActionChosen.

**run** (*n\_epochs*, *epoch\_length*)  
This function encapsulates the whole process of the learning. It starts by calling the controllers method “onStart”, Then it runs a given number of epochs where an epoch is made up of one or many episodes (called with agent.\_runEpisode) and where an epoch ends up after the number of steps reaches the argument “epoch\_length”. It ends up by calling the controllers method “end”.

**n\_epochs** [int] number of epochs

**epoch\_length** [int] maximum number of steps for a given epoch

**setControllersActive** (*toDisable*, *active*)  
Activate controller

**setDiscountFactor** (*df*)  
Set the discount factor

**setLearningRate** (*lr*)  
Set the learning rate for the gradient descent

**setNetwork** (*fname*, *nEpoch=-1*)  
Set values into the network

**fname** [string] Name of the file where the values are

**nEpoch** [int] Epoch number (Optional)

**totalRewardOverLastTest** ()  
Returns the average sum of rewards per episode and the number of episode

**train** ()  
This function selects a random batch of data (with self.\_dataset.randomBatch) and performs a Q-learning iteration (with self.\_learning\_algo.train).

---

```
class deer.agent.DataSet (env, random_state=None, max_size=1000000, use_priority=False,
                        only_full_history=True)
```

A replay memory consisting of circular buffers for observations, actions, rewards and terminals.

```
actions ()
```

Get all actions currently in the replay memory, ordered by time where they were taken.

```
addSample (obs, action, reward, is_terminal, priority)
```

Store the punctual observations, action, reward, is\_terminal and priority in the dataset. Parameters ———

**obs** : ndarray

An ndarray(dtype='object') where obs[s] corresponds to the punctual observation s before the agent took action [action].

**action** [int] The action taken after having observed [obs].

**reward** [float] The reward associated to taking this [action].

**is\_terminal** [bool] Tells whether [action] lead to a terminal state (i.e. corresponded to a terminal transition).

**priority** [float] The priority to be associated with the sample

```
observations ()
```

Get all observations currently in the replay memory, ordered by time where they were observed.

```
randomBatch (batch_size, use_priority)
```

Returns a batch of states, actions, rewards, terminal status, and next\_states for a number batch\_size of randomly chosen transitions. Note that if terminal[i] == True, then next\_states[s][i] == np.zeros\_like(states[s][i]) for each s.

**batch\_size** [int] Number of transitions to return.

**use\_priority** [Boolean] Whether to use prioritized replay or not

**states** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)]. States are taken randomly in the data with the only constraint that they are complete regarding the history size for each observation.

**actions** [numpy array of integers [batch\_size]] actions[i] is the action taken after having observed states[:,i].

**rewards** [numpy array of floats [batch\_size]] rewards[i] is the reward obtained for taking actions[i-1].

**next\_states** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**terminals** [numpy array of booleans [batch\_size] ] terminals[i] is True if the transition leads to a terminal state and False otherwise

**SliceError** If a batch of this batch\_size could not be built based on current data set (not enough data or all trajectories are too short).

```
randomBatch_nstep (batch_size, nstep, use_priority)
```

Return corresponding states, actions, rewards, terminal status, and next\_states for a number batch\_size of randomly chosen transitions. Note that if terminal[i] == True, then next\_states[s][i] == np.zeros\_like(states[s][i]) for each s.

**batch\_size** [int] Number of transitions to return.

**nstep** [int] Number of transitions to be considered for each element

**use\_priority** [Boolean] Whether to use prioritized replay or not

**states** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* (history size+nstep-1) \* size of punctual observation (which is 2D,1D or scalar)]. States are taken randomly in the data with the only constraint that they are complete regarding the history size for each observation.

**actions** [numpy array of integers [batch\_size, nstep]] actions[i] is the action taken after having observed states[:,i].

**rewards** [numpy array of floats [batch\_size, nstep]] rewards[i] is the reward obtained for taking actions[i-1].

**next\_states** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* (history size+nstep-1) \* size of punctual observation (which is 2D,1D or scalar)].

**terminals** [numpy array of booleans [batch\_size, nstep] ] terminals[i] is True if the transition leads to a terminal state and False otherwise

**SliceError** If a batch of this size could not be built based on current data set (not enough data or all trajectories are too short).

**rewards** ()

Get all rewards currently in the replay memory, ordered by time where they were received.

**terminals** ()

Get all terminals currently in the replay memory, ordered by time where they were observed.

terminals[i] is True if actions()[i] lead to a terminal state (i.e. corresponded to a terminal transition), and False otherwise.

**updatePriorities** (priorities, rndValidIndices)

## 3.2 Controller

This file defines the base Controller class and some presets controllers that you can use for controlling the training and the various parameters of your agents.

Controllers can be attached to an agent using the agent's `attach(Controller)` method. The order in which controllers are attached matters. Indeed, if controllers C1, C2 and C3 were attached in this order and C1 and C3 both listen to the `onEpisodeEnd` signal, the `onEpisodeEnd()` method of C1 will be called *before* the `onEpisodeEnd()` method of C3, whenever an episode ends.

<code>Controller()</code>	A base controller that does nothing when receiving the various signals emitted by an agent.
<code>LearningRateController(...)</code>	A controller that modifies the learning rate periodically upon epochs end.
<code>EpsilonController([initial_e, e_decays, ...])</code>	A controller that modifies the probability “epsilon” of taking a random action periodically.
<code>DiscountFactorController(...)</code>	A controller that modifies the q-network discount periodically.
<code>TrainerController([evaluate_on, ...])</code>	A controller that makes the agent train on its current database periodically.

Continued on next page

Table 2 – continued from previous page

<i>InterleavedTestEpochController</i> ([id, ...])	A controller that interleaves a test epoch between training epochs of the agent.
<i>FindBestController</i> ([validationID, testID, ...])	A controller that finds the neural net performing at best in validation mode (i.e.

**class** `deer.experiment.base_controllers.Controller`

A base controller that does nothing when receiving the various signals emitted by an agent. This class should be the base class of any controller you would want to define.

**onActionChosen** (*agent*, *action*)

Called whenever the agent has chosen an action.

This occurs after the agent state was updated with the new observation it made, but before it applied this action on the environment and before the total reward is updated.

**onActionTaken** (*agent*)

Called whenever the agent has taken an action on its environment.

This occurs after the agent applied this action on the environment and before terminality is evaluated. This is called only once, even in the case where the agent skip frames by taking the same action multiple times. In other words, this occurs just before the next observation of the environment.

**onEnd** (*agent*)

Called when the agent has finished processing all its epochs, just before returning from its `run()` method.

**onEpisodeEnd** (*agent*, *terminal\_reached*, *reward*)

Called whenever the agent ends an episode, just after this episode ended and before any `onEpochEnd()` signal could be sent.

**agent** [NeuralAgent] The agent firing the event

**terminal\_reached** [bool] Whether the episode ended because a terminal transition occurred. This could be False if the episode was stopped because its step budget was exhausted.

**reward** [float] The reward obtained on the last transition performed in this episode.

**onEpochEnd** (*agent*)

Called whenever the agent ends an epoch, just after the last episode of this epoch was ended and after any `onEpisodeEnd()` signal was processed.

**agent** [NeuralAgent] The agent firing the event

**onStart** (*agent*)

Called when the agent is going to start working (before anything else).

This corresponds to the moment where the agent's `run()` method is called.

**agent** [NeuralAgent] The agent firing the event

**setActive** (*active*)

Activate or deactivate this controller.

A controller should not react to any signal it receives as long as it is deactivated. For instance, if a controller maintains a counter on how many episodes it has seen, this counter should not be updated when this controller is disabled.

**class** `deer.experiment.base_controllers.LearningRateController` (*initial\_learning\_rate=0.005*,  
*learning\_rate\_decay=1.0*,  
*periodicity=1*)

Bases: `deer.experiment.base_controllers.Controller`

A controller that modifies the learning rate periodically upon epochs end.

**initial\_learning\_rate** [float] The learning rate upon agent start

**learning\_rate\_decay** [float] The factor by which the previous learning rate is multiplied every [periodicity] epochs.

**periodicity** [int] How many epochs are necessary before an update of the learning rate occurs

```
class deer.experiment.base_controllers.EpsilonController (initial_e=1.0,
                                                         e_decays=10000,
                                                         e_min=0.1,      eval-
                                                         uate_on='action',
                                                         periodicity=1,    re-
                                                         set_every='none')
```

Bases: `deer.experiment.base_controllers.Controller`

A controller that modifies the probability “epsilon” of taking a random action periodically.

**initial\_e** [float] Start epsilon

**e\_decays** [int] How many updates are necessary for epsilon to reach eMin

**e\_min** [float] End epsilon

**evaluate\_on** [str] After what type of event epsilon should be updated periodically. Possible values: ‘action’, ‘episode’, ‘epoch’.

**periodicity** [int] How many [evaluateOn] are necessary before an update of epsilon occurs

**reset\_every** [str] After what type of event epsilon should be reset to its initial value. Possible values: ‘none’, ‘episode’, ‘epoch’.

```
class deer.experiment.base_controllers.DiscountFactorController (initial_discount_factor=0.9,
                                                                    dis-
                                                                    count_factor_growth=1.0,
                                                                    dis-
                                                                    count_factor_max=0.99,
                                                                    periodic-
                                                                    ity=1)
```

Bases: `deer.experiment.base_controllers.Controller`

A controller that modifies the q-network discount periodically. More informations in : Francois-Lavet Vincent et al. (2015) - How to Discount Deep Reinforcement Learning: Towards New Dynamic Strategies (<http://arxiv.org/abs/1512.02011>).

**initial\_discount\_factor** [float] Start discount

**discount\_factor\_growth** [float] The factor by which the previous discount is multiplied every [periodicity] epochs.

**discount\_factor\_max** [float] Maximum reachable discount

**periodicity** [int] How many training epochs are necessary before an update of the discount occurs

```
class deer.experiment.base_controllers.TrainerController (evaluate_on='action',
                                                            periodicity=1,
                                                            show_episode_avg_V_value=True,
                                                            show_avg_Bellman_residual=True)
```

Bases: `deer.experiment.base_controllers.Controller`

A controller that makes the agent train on its current database periodically.



**evaluate\_on** [str] After what type of event the agent should be trained periodically. Possible values: ‘action’, ‘episode’, ‘epoch’. The first training will occur after the first occurrence of [evaluateOn].

**periodicity** [int] How many [evaluateOn] are necessary before a training occurs **\_show\_avg\_Bellman\_residual** [bool] - Whether to show an informative message after each episode end (and after a training if [evaluateOn] is ‘episode’) about the average bellman residual of this episode

**show\_episode\_avg\_V\_value** [bool] Whether to show an informative message after each episode end (and after a training if [evaluateOn] is ‘episode’) about the average V value of this episode

```
class deer.experiment.base_controllers.InterleavedTestEpochController (id=0,
                                                                    epoch_length=500,
                                                                    con-
                                                                    trollers_to_disable=[],
                                                                    peri-
                                                                    odic-
                                                                    ity=2,
                                                                    show_score=True,
                                                                    sum-
                                                                    ma-
                                                                    size_every=10)
```

Bases: *deer.experiment.base\_controllers.Controller*

A controller that interleaves a test epoch between training epochs of the agent.

**id** [int] The identifier ( $\geq 0$ ) of the mode each test epoch triggered by this controller will belong to. Can be used to discriminate between datasets in your Environment subclass (this is the argument that will be given to your environment’s reset() method when starting the test epoch).

**epoch\_length** [float] The total number of transitions that will occur during a test epoch. This means that this epoch could feature several episodes if a terminal transition is reached before this budget is exhausted.

**controllers\_to\_disable** [list of int] A list of controllers to disable when this controller wants to start a test epoch. These same controllers will be reactivated after this controller has finished dealing with its test epoch.

**periodicity** [int] How many epochs are necessary before a test epoch is ran (these controller’s epochs included: “1 test epoch on [periodicity] epochs”). Minimum value: 2.

**show\_score** [bool] Whether to print an informative message on stdout at the end of each test epoch, about the total reward obtained in the course of the test epoch.

**summarize\_every** [int] How many of this controller’s test epochs are necessary before the attached agent’s summarizeTestPerformance() method is called. Give a value  $\leq 0$  for “never”. If  $> 0$ , the first call will occur just after the first test epoch.

```
class deer.experiment.base_controllers.FindBestController (validationID=0,
                                                            testID=None,
                                                            unique_fname='nnet')
```

Bases: *deer.experiment.base\_controllers.Controller*

A controller that finds the neural net performing at best in validation mode (i.e. for mode = [validationID]) and computes the associated generalization score in test mode (i.e. for mode = [testID], and this only if [testID] is different from None). This controller should never be disabled by InterleavedTestControllers as it is meant to work in conjunction with them.

At each epoch end where this controller is active, it will look at the current mode the agent is in.

If the mode matches [validationID], it will take the total reward of the agent on this epoch and compare it to its current best score. If it is better, it will ask the agent to dump its current nnet on disk and update its current best score. In all cases, it saves the validation score obtained in a vector.

If the mode matches [testID], it saves the test (= generalization) score in another vector. Note that if [testID] is None, no test mode score are ever recorded.

At the end of the experiment (onEnd), if active, this controller will print information about the epoch at which the best neural net was found together with its generalization score, this last information shown only if [testID] is different from None. Finally it will dump a dictionary containing the data of the plots ({n: number of epochs elapsed, ts: test scores, vs: validation scores}). Note that if [testID] is None, the value dumped for the 'ts' key is [].

**validationID** [int] See synopsis

**testID** [int] See synopsis

**unique\_fname** [str] A unique filename (basename for score and network dumps).

### 3.3 Environment

This module defines the base class for the environments.

**class** deer.base\_classes.Environment

All your Environment classes should inherit this interface.

The environment defines the dynamics and the reward signal that the agent observes when interacting with it.

An agent sees at any time-step from the environment a collection of observable elements. Observing the environment at time t thus corresponds to obtaining a punctual observation for each of these elements. According to the control problem to solve, it might be useful for the agent to not only take action based on the current punctual observations but rather on a collection of the last punctual observations. In this framework, it's the environment that defines the number of each punctual observation to be considered.

Different "modes" are used in this framework to allow the environment to have different dynamics and/or reward signal. For instance, in training mode, only a part of the dynamics may be available so that it is possible to see how well the agent generalizes to a slightly different one.

**act** (action)

Applies the agent action [action] on the environment.

**action** [int] The action selected by the agent to operate on the environment. Should be an identifier included between 0 included and nActions() excluded.

**end** ()

Optional hook called at the end of all epochs

**inTerminalState** ()

Tells whether the environment reached a terminal state after the last transition (i.e. the last transition that occurred was terminal).

As the majority of control tasks considered have no end (a continuous control should be operated), by default this returns always False. But in the context of a video game for instance, terminal states can happen and in these cases, this method should be overridden.

**isTerminal** [bool] Whether or not the current state is terminal

**inputDimensions** ()

Gets the shape of the input space for this environment.

This returns a list whose length is the number of observations in the environment. Each element of the list is a tuple: the first integer is always the history size considered for this observation and the rest describes the shape of the observation at a given time step. For instance: - () or (1,) means each observation at a given time step is a single scalar and the history size is 1 (= only current observation) - (N,) means each

observation at a given time step is a single scalar and the history size is N - (N, M) means each observation at a given time step is a vector of length M and the history size is N - (N, M1, M2) means each observation at a given time step is a 2D matrix with M1 rows and M2 columns and the history size is N

#### **nActions()**

Gets the number of different actions that can be taken on this environment. It can be either an integer in the case of a finite discrete number of actions or it can be a list of couples [min\_action\_value,max\_action\_value] for a continuous action space

#### **observationType(subject)**

Gets the most inner type (np.uint8, np.float32, ...) of [subject].

**subject** [int] The subject

#### **observe()**

Gets a list of punctual observations composing this environment.

This returns a list where element i is a punctual observation. Note that the history of observations is not returned and only the current observation is.

See the documentation of inputDimensions() for more information about the shape of the observations.

#### **reset(mode)**

Resets the environment and put it in mode [mode]. This function is called when beginning every new episode.

The [mode] can be used to discriminate for instance between an agent which is training or trying to get a validation or generalization score. The mode the environment is in should always be redefined by resetting the environment using this method, meaning that the mode should be preserved until the next call to reset().

**mode** [int] The mode to put the environment into. Mode “-1” is reserved and always means “training”.

Initialization of the pseudo state at the beginning of a new episode: list (of lists) with size given by inputDimensions

#### **summarizePerformance(test\_data\_set, \*args, \*\*kwargs)**

Optional hook that can be used to show a summary of the performance of the agent on the environment in the current mode.

**test\_data\_set** [agent.DataSet] The dataset maintained by the agent in the current mode, which contains observations, actions taken and rewards obtained, as well as whether each transition was terminal or not. Refer to the documentation of agent.DataSet for more information.

## 3.4 Learning algorithms

<code>deer.base_classes. LearningAlgo(environment, ...)</code>	All the Q-networks, actor-critic networks, etc.
<code>deer.learning_algos.q_net_keras. MyQNetwork(...)</code>	Deep Q-learning network using Keras (with any back-end)
<code>deer.learning_algos.AC_net_keras. MyACNetwork(...)</code>	Actor-critic learning (using Keras) with Deep Deterministic Policy Gradient (DDPG) for the continuous action domain
<code>deer.learning_algos.CRAR_keras. CRAR(environment)</code>	Combined Reinforcement learning via Abstract Representations (CRAR) using Keras

**class** deer.base\_classes.**LearningAlgo** (environment, batch\_size)

All the Q-networks, actor-critic networks, etc. should inherit this interface.

**environment** [object from class Environment] The environment linked to the Q-network

**batch\_size** [int] Number of tuples taken into account for each iteration of gradient descent

**chooseBestAction** (*state*)

Get the best action for a pseudo-state

**discountFactor** ()

Getting the discount factor

**learningRate** ()

Getting the learning rate

**qValues** (*state*)

Get the q value for one pseudo-state

**setDiscountFactor** (*df*)

Setting the discount factor

**df** [float] The discount factor that has to be set

**setLearningRate** (*lr*)

Setting the learning rate NB: The learning rate has usually to be set in the optimizer, hence this function should be overridden. Otherwise, the learning rate change is likely not to be taken into account

**lr** [float] The learning rate that has to be set

**train** (*states, actions, rewards, nextStates, terminals*)

This method performs the training step (e.g. using Bellman iteration in a deep Q-network) for one batch of tuples.

```
class deer.learning_algos.q_net_keras.MyQNetwork (environment, rho=0.9,
                                                rms_epsilon=0.0001, momentum=0,
                                                clip_norm=0, freeze_interval=1000,
                                                batch_size=32, update_rule='rmsprop',
                                                random_state=<mtrand.RandomState
object>, double_Q=False,
                                                neural_network=<class
'deer.learning_algos.NN_keras.NN'>)
```

Deep Q-learning network using Keras (with any backend)

**environment** [object from class Environment] The environment in which the agent evolves.

**rho** [float] Parameter for rmsprop. Default : 0.9

**rms\_epsilon** [float] Parameter for rmsprop. Default : 0.0001

**momentum** [float] Momentum for SGD. Default : 0

**clip\_norm** [float] The gradient tensor will be clipped to a maximum L2 norm given by this value.

**freeze\_interval** [int] Period during which the target network is freezed and after which the target network is updated. Default : 1000

**batch\_size** [int] Number of tuples taken into account for each iteration of gradient descent. Default : 32

**update\_rule**: str {sgd,rmsprop}. Default : rmsprop

random\_state : numpy random number generator double\_Q : bool, optional

Activate or not the double\_Q learning. More informations in : Hado van Hasselt et al. (2015) - Deep Reinforcement Learning with Double Q-learning.

**neural\_network** [object, optional] default is deer.learning\_algos.NN\_keras

**chooseBestAction** (*state*, \**args*, \*\**kwargs*)

Get the best action for a pseudo-state

*state* : one pseudo-state

The best action : int

**getAllParams** ()

Get all parameters used by the learning algorithm

Values of the parameters: list of numpy arrays

**qValues** (*state\_val*)

Get the q values for one belief state

*state\_val* : one belief state

The q values for the provided belief state

**setAllParams** (*list\_of\_values*)

Set all parameters used by the learning algorithm

**list\_of\_values** [list of numpy arrays] list of the parameters to be set (same order than given by getAllParams()).

**train** (*states\_val*, *actions\_val*, *rewards\_val*, *next\_states\_val*, *terminals\_val*)

Train the Q-network from one batch of data.

**states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**actions\_val** [numpy array of integers with size [self.\_batch\_size]] actions[i] is the action taken after having observed states[:,i].

**rewards\_val** [numpy array of floats with size [self.\_batch\_size]] rewards[i] is the reward obtained for taking actions[i-1].

**next\_states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**terminals\_val** [numpy array of booleans with size [self.\_batch\_size]] terminals[i] is True if the transition leads to a terminal state and False otherwise

Average loss of the batch training (RMSE) Individual (square) losses for each tuple

```
class deer.learning_algos.AC_net_keras.MyACNetwork (environment, rho=0.9,
                                                    rms_epsilon=0.0001, mo-
                                                    mentum=0, clip_norm=0,
                                                    freeze_interval=1000,
                                                    batch_size=32, up-
                                                    date_rule='rmsprop', ran-
                                                    dom_state=<mtrand.RandomState
                                                    object>, double_Q=False,
                                                    neural_network_critic=<class
                                                    'deer.learning_algos.NN_keras.NN'>,
                                                    neural_network_actor=<class
                                                    'deer.learning_algos.NN_keras.NN'>)
```

Actor-critic learning (using Keras) with Deep Deterministic Policy Gradient (DDPG) for the continuous action domain

**environment** [object from class Environment] The environment in which the agent evolves.

**rho** [float] Parameter for rmsprop. Default : 0.9

**rms\_epsilon** [float] Parameter for rmsprop. Default : 0.0001

**momentum** [float] Momentum for SGD. Default : 0

**clip\_norm** [float] The gradient tensor will be clipped to a maximum L2 norm given by this value.

**freeze\_interval** [int] Period during which the target network is freezed and after which the target network is updated. Default : 1000

**batch\_size** [int] Number of tuples taken into account for each iteration of gradient descent. Default : 32

**update\_rule:** str {sgd,rmsprop}. Default : rmsprop

**random\_state** [numpy random number generator] Set the random seed.

**double\_Q** [bool, optional] Activate or not the double\_Q learning. More informations in : Hado van Hasselt et al. (2015) - Deep Reinforcement Learning with Double Q-learning.

**neural\_network\_critic** [object, optional] default is deer.learning\_algos.NN\_keras

**neural\_network\_actor** [object, optional] default is deer.learning\_algos.NN\_keras

**chooseBestAction** (*state, \*args, \*\*kwargs*)  
Get the best action for a pseudo-state  
  
state : one pseudo-state  
  
best\_action : float estim\_value : float

**clip\_action** (*action*)  
Clip the possible actions if it is outside the action space defined by self.\_nActions self.\_nActions is given as [[low\_action1,high\_action1],[low\_action2,high\_action2], ...]

**getAllParams** ()  
Get all parameters used by the learning algorithm  
  
Values of the parameters: list of numpy arrays

**gradients** (*states, actions*)  
Returns the gradients on the Q-network for the different actions (used for policy update)

**setAllParams** (*list\_of\_values*)  
Set all parameters used by the learning algorithm  
  
**list\_of\_values** [list of numpy arrays] list of the parameters to be set (same order than given by getAllParams()).

**train** (*states\_val, actions\_val, rewards\_val, next\_states\_val, terminals\_val*)  
Train the actor-critic algorithm from one batch of data.  
  
**states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].  
  
**actions\_val** [numpy array of integers with size [self.\_batch\_size]] actions[i] is the action taken after having observed states[:,i].  
  
**rewards\_val** [numpy array of floats with size [self.\_batch\_size]] rewards[i] is the reward obtained for taking actions[i-1].  
  
**next\_states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].  
  
**terminals\_val** [numpy array of booleans with size [self.\_batch\_size]] terminals[i] is True if the transition leads to a terminal state and False otherwise  
  
Average loss of the batch training Individual losses for each tuple

```
class deer.learning_algos.CRAR_keras.CRAR(environment, rho=0.9, rms_epsilon=0.0001,
                                         momentum=0, clip_norm=0,
                                         freeze_interval=1000, batch_size=32,
                                         update_rule='rmsprop', random_state=<mtrand.RandomState object>,
                                         double_Q=False, neural_network=<class
                                         'deer.learning_algos.NN_CRAR_keras.NN'>,
                                         **kwargs)
```

Combined Reinforcement learning via Abstract Representations (CRAR) using Keras

**environment** [object from class Environment] The environment in which the agent evolves.

**rho** [float] Parameter for rmsprop. Default : 0.9

**rms\_epsilon** [float] Parameter for rmsprop. Default : 0.0001

**momentum** [float] Momentum for SGD. Default : 0

**clip\_norm** [float] The gradient tensor will be clipped to a maximum L2 norm given by this value.

**freeze\_interval** [int] Period during which the target network is freezed and after which the target network is updated. Default : 1000

**batch\_size** [int] Number of tuples taken into account for each iteration of gradient descent. Default : 32

**update\_rule: str** {sgd,rmsprop}. Default : rmsprop

**random\_state** [numpy random number generator] Set the random seed.

**double\_Q** [bool, optional] Activate or not the double\_Q learning. More informations in : Hado van Hasselt et al. (2015) - Deep Reinforcement Learning with Double Q-learning.

**neural\_network** [object, optional] Default is deer.learning\_algos.NN\_keras

**chooseBestAction** (*state*, *mode*, *\*args*, **\*\*kwargs**)

Get the best action for a pseudo-state

**state** [list of numpy arrays] One pseudo-state. The number of arrays and their dimensions matches self.environment.inputDimensions().

**mode** [int] Identifier of the mode (-1 is reserved for the training mode).

The best action : int

**getAllParams** ()

Provides all parameters used by the learning algorithm

Values of the parameters: list of numpy arrays

**qValues** (*state\_val*)

Get the q values for one pseudo-state (without planning)

**state\_val** [array of objects (or list of objects)] Each object is a numpy array that relates to one of the observations with size [1 \* history size \* size of punctual observation (which is 2D,1D or scalar)].

The q values for the provided pseudo state

**qValues\_planning** (*state\_val*, *R*, *gamma*, *T*, *Q*, *d*=5)

Get the average Q-values up to planning depth d for one pseudo-state.

**state\_val** [array of objects (or list of objects)] Each object is a numpy array that relates to one of the observations with size [1 \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**R** [float\_model] Model that fits the reward

**gamma** [float\_model] Model that fits the discount factor

**T** [transition\_model] Model that fits the transition between abstract representation

**Q** [Q\_model] Model that fits the optimal Q-value

**d** [int] planning depth

The average q values with planning depth up to d for the provided pseudo-state

**qValues\_planning\_abstr** (*state\_abstr\_val, R, gamma, T, Q, d, branching\_factor=None*)

Get the q values for pseudo-state(s) with a planning depth d. This function is called recursively by decreasing the depth d at every step.

state\_abstr\_val : internal state(s). R : float\_model

Model that fits the reward

**gamma** [float\_model] Model that fits the discount factor

**T** [transition\_model] Model that fits the transition between abstract representation

**Q** [Q\_model] Model that fits the optimal Q-value

**d** [int] planning depth

The Q-values with planning depth d for the provided encoded state(s)

**setAllParams** (*list\_of\_values*)

Set all parameters used by the learning algorithm

**list\_of\_values** [list of numpy arrays] list of the parameters to be set (same order than given by getAllParams()).

**setLearningRate** (*lr*)

Setting the learning rate

**lr** [float] The learning rate that has to be set

**train** (*states\_val, actions\_val, rewards\_val, next\_states\_val, terminals\_val*)

Train CRAR from one batch of data.

**states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**actions\_val** [numpy array of integers with size [self.\_batch\_size]] actions[i] is the action taken after having observed states[:,i].

**rewards\_val** [numpy array of floats with size [self.\_batch\_size]] rewards[i] is the reward obtained for taking actions[i-1].

**next\_states\_val** [numpy array of objects] Each object is a numpy array that relates to one of the observations with size [batch\_size \* history size \* size of punctual observation (which is 2D,1D or scalar)].

**terminals\_val** [numpy array of booleans with size [self.\_batch\_size]] terminals[i] is True if the transition leads to a terminal state and False otherwise

Average loss of the batch training for the Q-values (RMSE) Individual (square) losses for the Q-values for each tuple

## 3.5 Policies



---

<code>deer.base_classes.Policy(learning_algo, ...)</code>	Abstract class for all policies.
<code>deer.policies.EpsilonGreedyPolicy(...)</code>	The policy acts greedily with probability $1 - \epsilon$ and acts randomly otherwise.
<code>deer.policies.LongerExplorationPolicy(...)</code>	Simple alternative to $\epsilon$ -greedy that can explore more efficiently for a broad class of realistic problems.

---

**class** `deer.base_classes.Policy` (*learning\_algo, n\_actions, random\_state*)

Abstract class for all policies. A policy takes observations as input, and outputs an action.

`learning_algo` : object from class `LearningALgo` `n_actions` : int or list

Definition of the action space provided by `Environment.nActions()`

`random_state` : numpy random number generator

**action** (*state*)

Main method of the Policy class. It can be called by `agent.py`, given a state, and should return a valid action w.r.t. the environment given to the constructor.

**bestAction** (*state, mode=None, \*args, \*\*kwargs*)

Returns the best Action for the given state. This is an additional encapsulation for q-network.

**randomAction** ()

Returns a random action

**class** `deer.policies.EpsilonGreedyPolicy` (*learning\_algo, n\_actions, random\_state, epsilon*)

Bases: `deer.base_classes.policy.Policy`

The policy acts greedily with probability  $1 - \epsilon$  and acts randomly otherwise. It is now used as a default policy for the neural agent.

**epsilon** [float] Proportion of random steps

**action** (*state, mode=None, \*args, \*\*kwargs*)

Main method of the Policy class. It can be called by `agent.py`, given a state, and should return a valid action w.r.t. the environment given to the constructor.

**epsilon** ()

Get the epsilon for  $\epsilon$ -greedy exploration

**setEpsilon** (*e*)

Set the epsilon used for  $\epsilon$ -greedy exploration

**class** `deer.policies.LongerExplorationPolicy` (*learning\_algo, n\_actions, random\_state, epsilon, length=10*)

Bases: `deer.base_classes.policy.Policy`

Simple alternative to  $\epsilon$ -greedy that can explore more efficiently for a broad class of realistic problems.

**epsilon** [float] Proportion of random steps

**length** [int] Length of the exploration sequences that will be considered

**action** (*state, mode=None, \*args, \*\*kwargs*)

Main method of the Policy class. It can be called by `agent.py`, given a state, and should return a valid action w.r.t. the environment given to the constructor.

**epsilon** ()

Get the epsilon

**setEpsilon** (*e*)  
Set the epsilon

## CHAPTER 4

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



### d

`deer.agent`, [15](#)

`deer.base_classes.environment`, [22](#)

`deer.experiment.base_controllers`, [18](#)



## A

act() (deer.base\_classes.Environment method), 22  
 action() (deer.base\_classes.Policy method), 29  
 action() (deer.policies.EpsilonGreedyPolicy method), 29  
 action() (deer.policies.LongerExplorationPolicy method), 29  
 actions() (deer.agent.DataSet method), 17  
 addSample() (deer.agent.DataSet method), 17  
 avgBellmanResidual() (deer.agent.NeuralAgent method), 16  
 avgEpisodeVValue() (deer.agent.NeuralAgent method), 16

## B

bestAction() (deer.base\_classes.Policy method), 29

## C

chooseBestAction() (deer.base\_classes.LearningAlgo method), 24  
 chooseBestAction() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26  
 chooseBestAction() (deer.learning\_algos.CRAR\_keras.CRAR method), 27  
 chooseBestAction() (deer.learning\_algos.q\_net\_keras.MyQNetwork method), 25  
 clip\_action() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26  
 Controller (class in deer.experiment.base\_controllers), 19  
 CRAR (class in deer.learning\_algos.CRAR\_keras), 26

## D

DataSet (class in deer.agent), 16  
 deer.agent (module), 15  
 deer.base\_classes.environment (module), 22  
 deer.experiment.base\_controllers (module), 18  
 discountFactor() (deer.agent.NeuralAgent method), 16  
 discountFactor() (deer.base\_classes.LearningAlgo method), 24

DiscountFactorController (class in deer.experiment.base\_controllers), 20  
 dumpNetwork() (deer.agent.NeuralAgent method), 16

## E

end() (deer.base\_classes.Environment method), 22  
 Environment (class in deer.base\_classes), 22  
 epsilon() (deer.policies.EpsilonGreedyPolicy method), 29  
 epsilon() (deer.policies.LongerExplorationPolicy method), 29  
 EpsilonController (class in deer.experiment.base\_controllers), 20  
 EpsilonGreedyPolicy (class in deer.policies), 29

## F

FindBestController (class in deer.experiment.base\_controllers), 21

## G

getACNetwork() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26  
 GetAllParams() (deer.learning\_algos.CRAR\_keras.CRAR method), 27  
 GetAllParams() (deer.learning\_algos.q\_net\_keras.MyQNetwork method), 25  
 gradients() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26

## I

inputDimensions() (deer.base\_classes.Environment method), 22  
 InterleavedTestEpochController (class in deer.experiment.base\_controllers), 21  
 inTerminalState() (deer.base\_classes.Environment method), 22

## L

LearningAlgo (class in deer.base\_classes), 23  
 learningRate() (deer.agent.NeuralAgent method), 16

learningRate() (deer.base\_classes.LearningAlgo method), 24  
 LearningRateController (class in deer.experiment.base\_controllers), 19  
 LongerExplorationPolicy (class in deer.policies), 29

## M

MyACNetwork (class in deer.learning\_algos.AC\_net\_keras), 25  
 MyQNetwork (class in deer.learning\_algos.q\_net\_keras), 24

## N

nActions() (deer.base\_classes.Environment method), 23  
 NeuralAgent (class in deer.agent), 15

## O

observations() (deer.agent.DataSet method), 17  
 observationType() (deer.base\_classes.Environment method), 23  
 observe() (deer.base\_classes.Environment method), 23  
 onActionChosen() (deer.experiment.base\_controllers.Controller method), 19  
 onActionTaken() (deer.experiment.base\_controllers.Controller method), 19  
 onEnd() (deer.experiment.base\_controllers.Controller method), 19  
 onEpisodeEnd() (deer.experiment.base\_controllers.Controller method), 19  
 onEpochEnd() (deer.experiment.base\_controllers.Controller method), 19  
 onStart() (deer.experiment.base\_controllers.Controller method), 19  
 overrideNextAction() (deer.agent.NeuralAgent method), 16

## P

Policy (class in deer.base\_classes), 29

## Q

qValues() (deer.base\_classes.LearningAlgo method), 24  
 qValues() (deer.learning\_algos.CRAR\_keras.CRAR method), 27  
 qValues() (deer.learning\_algos.q\_net\_keras.MyQNetwork method), 25  
 qValues\_planning() (deer.learning\_algos.CRAR\_keras.CRAR method), 27  
 qValues\_planning\_abstr() (deer.learning\_algos.CRAR\_keras.CRAR method), 28

## R

randomAction() (deer.base\_classes.Policy method), 29

randomBatch() (deer.agent.DataSet method), 17  
 randomBatch\_nstep() (deer.agent.DataSet method), 17  
 reset() (deer.base\_classes.Environment method), 23  
 rewards() (deer.agent.DataSet method), 18  
 run() (deer.agent.NeuralAgent method), 16

## S

setActive() (deer.experiment.base\_controllers.Controller method), 19  
 setAllParams() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26  
 setAllParams() (deer.learning\_algos.CRAR\_keras.CRAR method), 28  
 setAllParams() (deer.learning\_algos.q\_net\_keras.MyQNetwork method), 25  
 setControllersActive() (deer.agent.NeuralAgent method), 16  
 setDiscountFactor() (deer.agent.NeuralAgent method), 16  
 setDiscountFactor() (deer.base\_classes.LearningAlgo method), 24  
 setEpsilon() (deer.policies.EpsilonGreedyPolicy method), 29  
 setEpsilon() (deer.policies.LongerExplorationPolicy method), 29  
 setLearningRate() (deer.agent.NeuralAgent method), 16  
 setLearningRate() (deer.base\_classes.LearningAlgo method), 24  
 setLearningRate() (deer.learning\_algos.CRAR\_keras.CRAR method), 28  
 setNetwork() (deer.agent.NeuralAgent method), 16  
 summarizePerformance() (deer.base\_classes.Environment method), 23

## T

terminals() (deer.agent.DataSet method), 18  
 totalRewardOverLastTest() (deer.agent.NeuralAgent method), 16  
 train() (deer.agent.NeuralAgent method), 16  
 train() (deer.base\_classes.LearningAlgo method), 24  
 train() (deer.learning\_algos.AC\_net\_keras.MyACNetwork method), 26  
 train() (deer.learning\_algos.CRAR\_keras.CRAR method), 28  
 train() (deer.learning\_algos.q\_net\_keras.MyQNetwork method), 25  
 TrainerController (class in deer.experiment.base\_controllers), 20

## U

updatePriorities() (deer.agent.DataSet method), 18