

---

**GenRL**  
*Release 0.1*

**Mar 12, 2022**



---

# User Guide

---

<b>1 Features</b>	<b>1</b>
<b>2 Contents</b>	<b>3</b>
2.1 Installation . . . . .	3
2.2 About . . . . .	3
2.3 Tutorials . . . . .	5
2.4 Agents . . . . .	30
2.5 Environments . . . . .	73
2.6 Core . . . . .	82
2.7 Utilities . . . . .	97
2.8 Trainers . . . . .	101
2.9 Common . . . . .	106
<b>Python Module Index</b>	<b>113</b>
<b>Index</b>	<b>115</b>



# CHAPTER 1

---

## Features

---

- Unified Trainer and Logging class: code reusability and high-level UI
- Ready-made algorithm implementations: ready-made implementations of popular RL algorithms.
- Extensive Benchmarking
- Environment implementations
- Heavy Encapsulation useful for new algorithms



# CHAPTER 2

---

## Contents

---

## 2.1 Installation

### 2.1.1 PyPI Package

GenRL is compatible with Python 3.6 or later and also depends on `pytorch` and `openai-gym`. The easiest way to install GenRL is with pip, Python's preferred package installer.

```
$ pip install genrl
```

Note that GenRL is an active project and routinely publishes new releases. In order to upgrade GenRL to the latest version, use pip as follows.

```
$ pip install -U genrl
```

### 2.1.2 From Source

If you intend to install the latest unreleased version of the library (i.e from source), you can simply do:

```
$ git clone https://github.com/SforAiD1/genrl.git
$ cd genrl
$ python setup.py install
```

## 2.2 About

### 2.2.1 Introduction

Reinforcement Learning has taken massive leaps forward in extending current AI research. David Silver's paper on playing Atari with Deep Reinforcement Learning can be considered one of the seminal papers in establishing

a completely new landscape of Reinforcement Learning Research. With applications in Robotics, Healthcare and numerous other domains, RL has become the prime mechanism of modelling Sequential Decision Making through AI.

Yet, current libraries and resources in Reinforcement Learning are either very limited, messy and/or are scattered. OpenAI's Spinning Up is a great resource for getting started with Deep Reinforcement Learning but it fails to cover more basic concepts in Reinforcement Learning for e.g. Multi Armed Bandits. garage is a great resource for reproducing and evaluating RL algorithms but it fails to introduce a newbie to RL.

With GenRL, our goal is three-fold: - To educate the user about Reinforcement learning. - Easy to understand implementations of State of the Art Reinforcement Learning Algorithms. - Providing utilities for developing and evaluating new RL algorithms. Or in a sense be able to implement any new RL algorithm in less than 200 lines.

## 2.2.2 Policies and Values

Modern research on Reinforcement Learning is majorly based on Markov Decision Processes. Policy and Value Functions are one of the core parts of such a problem formulation. And so, policies and values form one of the core parts of our library.

## 2.2.3 Trainers and Loggers

### Trainers

Most current algorithms follow a standard procedure of training. Considering a classification between On-Policy and Off-Policy Algorithms, we provide high level APIs through Trainers which can be coupled with Agents and Environments for training seamlessly.

Lets take the example of an On-Policy Algorithm, Proximal Policy Optimization. In our Agent, we make sure to define three methods: `collect_rollouts`, `get_traj_loss` and finally `update_policy`.

The `OnPolicyTrainer` simply calls these functions and enables high level usage by simple defining of three methods.

### Loggers

At the moment, we support three different types of Loggers. `HumanOutputFormat`, `TensorboardLogger` and `CSVLogger`. Any of these loggers can be initialized really easily by the top level `Logger` class and specifying the individual formats in which logging should performed.

```
logger = Logger(logdir='logs/', formats=['stdout', 'tensorboard'])
```

After which logger can perform logging easily by providing it with dictionaries of data. For e.g.

```
logger.write({'logger':0})
```

Note: The Tensorboard logger requires an extra x-axis parameter, as it plots data rather than just show it in a tabular format.

## 2.2.4 Agent Encapsulation

WIP

## 2.2.5 Environments

Wrappers

## 2.3 Tutorials

### 2.3.1 Bandit Tutorials

#### Multi Armed Bandit Overview

##### Training an EpsilonGreedy agent on a Bernoulli Multi Armed Bandit

Multi armed bandits is one of the most basic problems in RL. Think of it like this, you have ‘n’ levers in front of you and each of these levers will give you a different reward. For the purposes of formalising the problem the reward is written down in terms of a reward function i.e., the probability of getting a reward when a lever is pulled.

Suppose you try out one of the levers and get a positive reward. What do you do next? Should you just keep pulling that lever every time or think what if there might be a better reward to pulling one of the other levers? This is the exploration - exploitation dilemma.

*Exploitation* - Utilise the information you have gathered till now, to make the best decision. In this case, after 1 try you know a lever is giving you a positive reward and you just *exploit* it further. Since you do not care about other arms if you keep *exploiting*, it is known as the greedy action.

*Exploration* - You explore the untried levers in an attempt to maybe discover another one which has a higher payout than the one you currently have some knowledge about. This is exploring all your options without worrying about the short-term rewards, in hope of finding a lever with a bigger reward, in the long run.

You have to use an algorithm which correctly trades off exploration and exploitation as we do not want a ‘greedy’ algorithm which only exploits and does not explore at all, because there are very high chances that it will converge to a sub-optimal policy. We do not want an algorithm that keeps exploring either as this would lead to sub-optimal rewards inspite of knowing the best action to be taken. In this case, the optimal policy will be to always pull the lever with the highest reward, but at the beginning we do not know the probability distribution of the rewards.

So, we want a policy which explores actively at the beginning, building up an estimate for the reward values(defined as *quality*) of all the actions, and then exploiting that from that time onwards.

A Bernoulli Multi-Armed Bandit has multiple arms with each having a different bernoulli distribution over its reward. Basically each arm has a probabiltiy associated with it which is the probability of getting a reward if that arm is pulled. Our aim is to find the arm which has the highest probabiltiy, thus giving us the maximum return.

Notation:

$Q_t(a)$ : Estimated quality of action ‘a’ at timestep ‘t’.

$q(a)$ : True value of action ‘a’.

We want our estimate  $Q_t(a)$  to be as close to the true value  $q(a)$  as possible, so we can make the correct decision.

Let the action with the maximum quality be  $a^*$ :

$$q^* = q(a^*)$$

Our goal is to find this  $q^*$ .

The ‘regret function’ is defined as the sum of ‘regret’ accumulated over all timesteps. This regret is the cost of not choosing the optimal arm and instead of exploring. Mathematically it can be written as:

$$L = E\left[\sum_{t=0}^T q^* - Q_t(a)\right]$$

Some policies which are effective at exploring are: 1. [Epsilon Greedy](#) 2. [Gradient Algorithm](#) 3. [UCB\(Upper Confidence Bound\)](#) 4. [Bayesian](#) 5. [Thompson Sampling](#)

Epsilon Greedy is the most basic exploratory policy which follows a simple principle to balance exploration and exploitation. It ‘exploits’ the current knowledge of the bandit most of the times, i.e. takes the action with the largest  $q$  value. But with a small probability  $\epsilon$ , it also explores a random action. The value of  $\epsilon$  signifies how much you want the agent explore. Higher the value, the more it explores. But remember you do not want an agent to explore too much even after it has a pretty confident estimate of the reward function, so the value of  $\epsilon$  should neither be too high nor too low!

For the bandit, you can set the number of bandits, number of arms, and also reward probabilities of each of these arms separately.

Code to train an Epsilon Greedy agent on a Bernoulli Multi-Armed Bandit:

```
import gym
import numpy as np

from genrl.bandit import BernoulliMAB, EpsGreedyMABA, MABTrainer

reward_probs = np.random.random(size=(bandits, arms))
bandit = BernoulliMAB(arms=5, reward_probs=reward_probs, context_type="int")
agent = EpsGreedyMABA(bandit, eps=0.05)

trainer = MABTrainer(agent, bandit)
trainer.train(timesteps=10000)
```

More details can be found in the docs for [BernoulliMAB](#), [EpsGreedyMABA](#), [MABTrainer](#).

You can also refer to the book “Reinforcement Learning: An Introduction”, Chapter 2 for further information on bandits.

## Contextual Bandits Overview

### Problem Setting

To get some background on the basic multi armed bandit problem, we recommend that you go through the [Multi Armed Bandit Overview](#) first. The contextual bandit (CB) problem varies from the basic case in that at each timestep, a context vector  $x \in \mathbb{R}^d$  is presented to the agent. The agent must then decide on an action  $a \in \mathcal{A}$  to take based on that context. After the action is taken, the reward  $r \in \mathbb{R}$  for only that action is revealed to the agent (a feature of all reinforcement learning problems). The aim of the agent remains the same - minimising regret and thus finding an optimal policy.

Here you still have the problem of exploration vs exploitation, but the agent also needs to find some relation between the context and reward.

### A Simple Example

Lets consider the simplest case of the CB problem. Instead of having only one  $k$ -armed bandit that needs to be solved, say we have  $m$  different  $k$ -armed Bernoulli bandits. At each timestep, the context presented is the number of the

bandit for which an action needs to be selected:  $i \in \mathbb{I}$  where  $0 < i \leq m$

Although real life CB problems usually have much higher dimensional contexts, such a toy problem can be useful for testing and debugging agents.

To instantiate a Bernoulli bandit with  $m = 10$  and  $k = 5$  (10 different 5-armed bandits) -

```
from genrl.bandit import BernoulliMAB

bandit = BernoulliMAB(bandits=10, arms=5, context_type="int")
```

Note that this is using the same BernoulliMAB as in the simple bandit case except that instead of the `bandits` argument defaulting to 1, we are explicitly saying we want multiple bandits (a contextual case)

Suppose you want to solve this bandit with a UCB based policy.

```
from genrl.bandit import UCBMABAgent

agent = UCBMABAgent(bandit)
context = bandit.reset()

action = agent.select_action(context)
new_context, reward = bandit.step(action)
```

To train the agent, you can set up a loop which calls the `update_params` method on the agent whenever you want to agent to learn from actions it has taken. For convenience it is highly recommended to use the `MABTrainer` in such cases.

## Data based Contextual Bandits

Lets consider a more realistic class of CB problem. In real life, you the CB setting is usually used to model recommendation or classification problems. Here, instead of getting an integer as the context, you will get a  $d$ -dimensional feature vector  $\mathbf{x} \in \mathbb{R}^d$ . This is also different from regular classification since you only get the reward  $r \in \mathbb{R}$  for the action you have taken.

While tabular solutions can work well for integer contexts (see the implementation of any `genrl.bandit.MABAgent` for details), when you have a high dimensional vector, the agent should be able to infer the complex relation between the contexts and rewards. This can be done by modelling a conditional distribution over rewards for each action given the context.

$$P(r|a, \mathbf{x})$$

There are many ways to do this. For a detailed explanation and comparison of contextual bandit methods you can refer to [this paper](#).

The following are the agents implemented in `genrl`

- Linear Posterior Inference
- Neural Network based Linear
- Variational
- Neural Network based Epsilon Greedy
- Bootstrap
- Parameter noise Sampling

You can find the tutorials for most of these in *Bandit Tutorials*.

All the methods which use neural networks, provide an option to train and evaluate with dropout, have a decaying learning rate and a limit for gradient clipping. The sizes of hidden layers for the networks can also be specified. Refer to docs of the specific agents to see how to use these options.

Individual agents will have other method specific parameters to control behavior. Although default values have been provided, it may be necessary to tune these for individual use cases.

The following bandits based on datasets are implemented in genrl

- Adult Census Income Dataset
- US Census Dataset
- Forest covertype Dataset
- MAGIC Gamma Telescope dataset
- Mushroom Dataset
- Statlog Space Shuttle Dataset

For each bandit, while instantiating an object you can either specify a path to the data file or pass download=True as an argument to download the data directly.

## Data based Bandit Example

For this example, we'll model the Statlog dataset as a bandit problem. You can read more about the bandit in the Statlog docs. In brief we have the number of arms as  $k = 7$  and dimension of context vector as  $d = 9$ . The agent will get a reward  $r = 1$  if it selects the correct arm else  $r = 0$ .

```
from genrl.bandit import StatlogDataBandit

bandit = StatlogDataBandit(download=True)
context = bandit.reset()
```

Suppose you want to solve this bandit with a Greedy neural network based policy.

```
from genrl.bandit import NeuralLinearPosteriorAgent

agent = NeuralLinearPosteriorAgent(bandit)
context = bandit.reset()

action = agent.select_action(context)
new_context, reward = bandit.step(action)
```

To train the agent, we highly recommend using the DCBTrainer. You can refer to the implementation of the train function to get an idea of how to implement your own training loop.

```
from genrl.bandit import DCBTrainer

trainer = DCBTrainer(agent, bandit)
trainer.train(timesteps=5000, batch_size=32)
```

## Further material about bandits

1. Deep Contextual Multi-armed Bandits, Collier and Llorens, 2018

2. Deep Bayesian Bandits Showdown, Riquelme et al, 2018
3. A Contextual Bandit Bake-off, Bietti et al, 2020

## UCB

### Training a UCB algorithm on a Bernoulli Multi-Armed Bandit

For an introduction to Multi Armed Bandits, refer to [Multi Armed Bandit Overview](#)

The UCB algorithm follows a basic principle - ‘Optimism in the face of uncertainty’. What this means is that we should always select the action whose reward we are most uncertain of. We quantify the uncertainty of taking an action by calculating an upper bound of the quality(reward) for that action. We then select the greedy action with respect to this upper bound.

Hoeffding’s inequality:

$$P[q(a) > Q_t(a) + U_t(a)] \leq e^{-2N_t(a)U_t(a)^2}$$

,

$q(a)$  is the quality of that action,

$Q_t(a)$  is the estimate of the quality of action ‘a’ at time ‘t’,

$U_t(a)$  is the upper bound for uncertainty for that action at time ‘t’,

$N_t(a)$  is the number of times action ‘a’ has been selected

$$e^{-2N_t(a)U_t(a)^2} = t^{-4}$$

$$U_t(a) = \sqrt{\frac{2\log t}{N_t(a)}}$$

Action taken:  $a = \text{argmax}(Q_t(a) + U_t(a))$

As we can see, the less an action has been tried, more the uncertainty is (due to  $N_t(a)$  being in the denominator), which leads to that action having a higher chance of being explored. Also, theoretically, as  $N_t(a)$  goes to infinity, the uncertainty decreases to 0 giving us the true value of the quality of that action:  $q(a)$ . This allows us to ‘exploit’ the greedy action  $a^*$  from then.

Code to train a UCB agent on a Bernoulli Multi-Armed Bandit:

```
import gym
import numpy as np

from genrl.bandit import BernoulliMAB, MABTrainer, UCBMABAgent

bandits = 10
arms = 5

reward_probs = np.random.random(size=(bandits, arms))
bandit = BernoulliMAB(bandits, arms, reward_probs, context_type="int")
agent = UCBMABAgent(bandit, confidence=1.0)

trainer = MABTrainer(agent, bandit)
trainer.train(timesteps=10000)
```

More details can be found in the docs for [BernoulliMAB](#), [UCB](#) and [MABTrainer](#).

## Thompson Sampling

### Using Thompson Sampling on a Bernoulli Multi-Armed Bandit

For an introduction to Multi Armed Bandits, refer to [Multi Armed Bandit Overview](#)

Thompson Sampling is one of the best methods for solving the Bernoulli multi-armed bandits problem. It is a ‘sample-based probability matching’ method.

We initially *assume* an initial distribution(prior) over the quality of each of the arms. We can model this prior using a Beta distribution, parametrised by alpha( $\alpha$ ) and beta( $\beta$ ).

$$PDF = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$$

Let’s just think of the denominator as some normalising constant, and focus on the numerator for now. We initialise  $\alpha = \beta = 1$ . This will result in a uniform distribution over the values (0, 1), making all the values of quality from 0 to 1 equally probable, so this is a fair initial assumption. Now think of  $\alpha$  as the number of times we get the reward ‘1’ and  $\beta$  as the number of times we get ‘0’, for a particular arm. As our agent interacts with the environment and gets a reward for pulling any arm, we will update our prior for that arm using Bayes Theorem. What this does is that it gives a posterior distribution over the quality, according to the rewards we have seen so far.

At each timestep, we sample the quality:  $Q_t(a)$  for each arm from the posterior and select the sample with the highest value. The more an action is tried out, the narrower is the distribution over its quality, meaning we have a confident estimate of its quality ( $q(a)$ ). If an action has not been tried out that often, it will have a more wider distribution (high variance), meaning we are uncertain about our estimate of its quality ( $q(a)$ ). This wider variance of an arm with an uncertain estimate creates opportunities for it to be picked during sampling.

As we experience more successes for a particular arm, the value of  $\alpha$  for that arm increases and similarly, the more failures we experience, the value of  $\beta$  increases. Higher the value of one of the parameters as compared to the other, the more skewed is the distribution in one of the directions. For eg. if  $\alpha = 100$  and  $\beta = 50$ , we have seen considerably more successes than failures for this arm and so our estimate for its quality should be  $>0.5$ . This will be reflected in the posterior of this arm, i.e. the mean of the distribution, characterised by  $\frac{\alpha}{\alpha+\beta}$  will be 0.66, which is  $>0.5$  as we expected.

Code to use Thompson Sampling on a Bernoulli Multi-Armed Bandit:

```
import gym
import numpy as np

from genrl.bandit import BernoulliMAB, MABTrainer, ThompsonSamplingMABAgent

bandits = 10
arms = 5
alpha = 1.0
beta = 1.0

reward_probs = np.random.random(size=(bandits, arms))
bandit = BernoulliMAB(bandits, arms, reward_probs, context_type="int")
agent = ThompsonSamplingMABAgent(bandit, alpha, beta)

trainer = MABTrainer(agent, bandit)
trainer.train(timesteps=10000)
```

More details can be found in the docs for [BernoulliMAB](#), [UCB](#) and [MABTrainer](#).

## Bayesian

## Using Bayesian Method on a Bernoulli Multi-Armed Bandit

For an introduction to Multi Armed Bandits, refer to [Multi Armed Bandit Overview](#)

This method is also based on the principle - ‘Optimism in the face of uncertainty’, like UCB. We initially *assume* an initial distribution(prior) over the quality of each of the arms. We can model this prior using a Beta distribution, parametrised by alpha( $\alpha$ ) and beta( $\beta$ ).

$$PDF = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$$

Let’s just think of the denominator as some normalising constant, and focus on the numerator for now. We initialise  $\alpha = \beta = 1$ . This will result in a uniform distribution over the values (0, 1), making all the values of quality from 0 to 1 equally probable, so this is a fair initial assumption. Now think of  $\alpha$  as the number of times we get the reward ‘1’ and  $\beta$  as the number of times we get ‘0’, for a particular arm. As our agent interacts with the environment and gets a reward for pulling any arm, we will update our prior for that arm using Bayes Theorem. What this does is that it gives a posterior distribution over the quality, according to the rewards we have seen so far.

This is quite similar to [Thompson Sampling](#). But what is different here is that we explicitly try to calculate the uncertainty of a particular action by calculating the standard deviation( $\sigma$ ) of its posterior. We add this std. dev to the mean of the posterior, giving us an *upper bound* of the quality of that arm. At each timestep we select a greedy action based on this upper bound we calculated.

$$a_t = argmax(q_t(a) + \sigma_{q_t})$$

As we try out an action more and more, the standard deviation of the posterior decreases, corresponding to a decrease in the uncertainty of that action, which is exactly what we want. If an action has not been tried that often, it will have a wider posterior, meaning higher chances of it getting selected based on its upper bound.

Code to use Bayesian method on a Bernoulli Multi-Armed Bandit:

```
import gym
import numpy as np

from genrl.bandit import BayesianUCBMABAgent, BernoulliMAB, MABTrainer

bandits = 10
arms = 5
alpha = 1.0
beta = 1.0

reward_probs = np.random.random(size=(bandits, arms))
bandit = BernoulliMAB(bandits, arms, reward_probs, context_type="int")
agent = BayesianUCBMABAgent(bandit, alpha, beta)

trainer = MABTrainer(agent, bandit)
trainer.train(timesteps=10000)
```

More details can be found in the docs for [BernoulliMAB](#), [BayesianUCBMABAgent](#) and [MABTrainer](#).

## Gradients

### Using Gradient Method on a Bernoulli Multi-Armed Bandit

For an introduction to Multi Armed Bandits, refer to [Multi Armed Bandit Overview](#)

This method is different compared to others. In other methods, we explicitly attempt to estimate the ‘value’ of taking an action (its quality) whereas in this method we approach the problem in a different way. Here, instead of estimating how

good an action is through its quality, we only care about its preference of being selected compared to other actions. We denote this preference by  $H_t(a)$ . The larger the preference of an action ‘a’, more are the chances of it being selected, but this preference has no interpretation in terms of the reward for that action. Only the relative preference is important.

The action probabilities are related to these action preferences  $H_t(a)$  by a softmax function. The probability of taking action  $a_j$  is given by:

$$P(a_j) = \frac{e^{H_t(a_j)}}{\sum_{i=1}^A e^{H_t(a_i)}} = \pi_t(a_j)$$

where, A is the total number of actions and  $\pi_t(a)$  is the probability of taking action ‘a’ at timestep ‘t’.

We initialise the preferences for all the actions to be 0, meaning  $\pi_t(a) = \frac{1}{A}$  for all actions.

After computing  $\pi_t(a)$  for all actions at each timestep, the action is sampled using this probability. Then that action is performed and based on the reward we get, we update our preferences.

The update rule basically performs stochastic gradient ascent:

$$H_{t+1}(a_t) = H_t(a_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(a_t)), \text{ for } a_t: \text{action taken at time 't'}$$

$$H_{t+1}(a) = H_t(a) - \alpha(R_t - \bar{R}_t)(\pi_t(a)) \text{ for rest of the actions}$$

where,  $\alpha$  is the step size,  $R_t$  is the reward obtained at time ‘t’ and  $\bar{R}_t$  is the mean reward obtained upto time t. If current reward is larger than the mean reward, we increase our preference for that action taken at time ‘t’. If it is lower than the mean reward, we decrease our preference for that action. The preferences for the rest of the actions are updated in the opposite direction.

For a more detailed mathematical analysis and derivation of the update rule, refer to chapter 2 of Sutton & Barto.

Code to use the Gradient method on a Bernoulli Multi-Armed Bandit:

```
import gym
import numpy as np

from genrl.bandit import BernoulliMAB, GradientMABAgent, MABTrainer

bandits = 10
arms = 5

reward_probs = np.random.random(size=(bandits, arms))
bandit = BernoulliMAB(bandits, arms, reward_probs, context_type="int")
agent = GradientMABAgent(bandit, alpha=0.1, temp=0.01)

trainer = MABTrainer(agent, bandit)
trainer.train(timesteps=10000)
```

More details can be found in the docs for [BernoulliMAB](#), [BayesianUCBMABAgent](#) and [MABTrainer](#).

## Linear Posterior Inference

For an introduction to the Contextual Bandit problem, refer to [Contextual Bandits Overview](#).

In this agent we assume a linear relationship between context and reward distribution of the form

$$Y = X^T \beta + \epsilon \text{ where } \epsilon \sim \mathcal{N}(0, \sigma^2)$$

We can utilise [bayesian linear regression](#) to find the parameters  $\beta$  and  $\sigma$ . Since our agent is continually learning, the parameters of the model will be updated according to the  $(x, a, r)$  transitions it observes.

For more complex non linear relations, we can make use of neural networks to transform the context into a learned embedding space. The above method can then be used on this latent embedding to model the reward.

An example of using a neural network based linear posterior agent in genrl -

```
from genrl.bandit import NeuralLinearPosteriorAgent, DCBTrainer

agent = NeuralLinearPosteriorAgent(bandit, lambda_prior=0.5, a0=2, b0=2, device="cuda")
trainer = DCBTrainer(agent, bandit)
trainer.train()
```

Note that the priors here are used to parameterise the initial distribution over  $\beta$  and  $\sigma$ . More specifically `lambda_prior` is used to parameterise a gaussian distribution for  $\beta$  while `a0` and `b0` are parameters of an inverse gamma distribution over  $\sigma^2$ . These are updated over the course of exploring a bandit. More details can be found in Section 3 of [this paper](#).

All hyperparameters can be tuned for individual use cases to improve training efficiency and achieve convergence faster.

Refer to the [LinearPosteriorAgent](#), [NeuralLinearPosteriorAgent](#) and [DCBTrainer](#) docs for more details.

## Variational Inference

For an introduction to the Contextual Bandit problem, refer to [Contextual Bandits Overview](#).

In this method, we try find a distribution  $P_\theta(r|x, a)$  by minimising the KL divergence with the true distribution. For the model we take a neural network where each weight is modelled by independent gaussians, also known as Bayesian Neural Nets.

An example of using a variational inference based agent in genrl with bayesian net of hidden layer of 128 neurons and standard deviation of 0.1 for all the weights -

```
from genrl.bandit import VariationalAgent, DCBTrainer

agent = VariationalAgent(bandit, hidden_dims=[128], noise_std=0.1, device="cuda")

trainer = DCBTrainer(agent, bandit)
trainer.train()
```

Refer to the [VariationalAgent](#), and [DCBTrainer](#) docs for more details.

## Bootstrap

For an introduction to the Contextual Bandit problem, refer to [Contextual Bandits Overview](#).

In the bootstrap agent multiple different neural network based models are trained simultaneously. Different transition databases are maintained for each model and every time we observe a transition it is added to each dataset with some probability. At each timestep, the model used to select an action is chosen randomly from the set of models.

By having multiple different models initialised with different random weights, we promote the exploration of the loss landscape which may have multiple different local optima.

An example of using a bootstrap based agent in genrl with 10 models with a hidden layer of 128 neurons which also uses dropout for training -

```
from genrl.bandit import BootstrapNeuralAgent, DCBTrainer

agent = BootstrapNeuralAgent(bandit, hidden_dims=[128], n=10, dropout_p=0.5, device=
    "cuda")

trainer = DCBTrainer(agent, bandit)
trainer.train()
```

Refer to the [BootstrapNeuralAgent](#) and [DCBTrainer](#) docs for more details.

## Parameter Noise Sampling

For an introduction to the Contextual Bandit problem, refer to [Contextual Bandits Overview](#).

One of the ways to improve exploration of our algorithms is to introduce noise into the weights of the neural network while selecting actions. This does not affect the gradients but will have a similar effect to epsilon greedy exploration.

The noise distribution is regularly updated during training to keep the KL divergence of the prediction and noise predictions within certain limits.

An example of using a noise sampling based agent in genrl with noise standard deviation as 0.1, KL divergence limit as 0.1 and batch size for updating the noise distribution as 128 -

```
from genrl.bandit import BootstrapNeuralAgent, DCBTrainer

agent = NeuralNoiseSamplingAgent(bandit, hidden_dims=[128], noise_std_dev=0.1, eps=0.
    1, noise_update_batch_size=128, device="cuda")

trainer = DCBTrainer(agent, bandit)
trainer.train()
```

Refer to the [NeuralNoiseSamplingAgent](#), and [DCBTrainer](#) docs for more details.

## Adding a new Data Bandit

The bandit submodule like all of genrl has been designed to be easily extensible for custom additions. This tutorial will show how to create a dataset based bandit which will work with the rest of genrl.bandit

For this tutorial, we will use the [Wine dataset](#) which is a simple dataset often used for testing classifiers. It has 178 examples each with 14 features, the first of which gives the cultivar of the wine (the feature we need to classify each wine sample into) (this can be one of three) and the rest give the properties of the wine itself. Formulated as a bandit problem we have a bandit with 3 arms and a 13-dimensional context. The agent will get a reward of 1 if it correctly selects the arm else 0.

To start off with lets import necessary modules, specify the data URL and make a class which inherits from `genrl.utils.data_bandits.base.DataBasedBandit`

```
from typing import Tuple

import pandas as pd
import torch

from genrl.utils.data_bandits.base import DataBasedBandit
from genrl.utils.data_bandits.utils import download_data
```

(continues on next page)

(continued from previous page)

```
URL = "http://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data"

class WineDataBandit(DataBasedBandit):
    def __init__(self, **kwargs):

        def reset(self) -> torch.Tensor:

            def _compute_reward(self, action: int) -> Tuple[int, int]:

                def _get_context(self) -> torch.Tensor:
```

We will need to implement `__init__`, `reset`, `_compute_reward` and `_get_context` to make the class functional.

For dataset based bandits, we can generally load the data into memory during initialisation. This can be in some tabular form (`numpy.array`, `torch.Tensor` or `pandas.DataFrame`) and maintaining an index. When reset, the bandit would set its index to 0 and reshuffle the rows of the table. For stepping, the bandit can compute rewards from the current row of the table as given by the index and then increment the index to move to the next row.

Lets start with `__init__`. Here we need to download the data if specified and load it into memory. Many utility functions are available in `genrl.utils.data_bandits.utils` including `download_data` to download data from a URL as well as functions to fetch data from memory.

For most cases, you can load the data into a `pandas.DataFrame`. You also need to specify the `n_actions`, `context_dim` and `len` here.

```
def __init__(self, **kwargs):
    super(WineDataBandit, self).__init__(**kwargs)

    path = kwargs.get("path", "./data/Wine/")
    download = kwargs.get("download", None)
    force_download = kwargs.get("force_download", None)
    url = kwargs.get("url", URL)

    if download:
        path = download_data(path, url, force_download)

    self._df = pd.read_csv(path, header=None)
    self.n_actions = len(self._df[0].unique())
    self.context_dim = self._df.shape[1] - 1
    self.len = len(self._df)
```

The `reset` method will shuffle the indices of the data and return the counting index to 0. You must have a call to `_reset` here to reset any metrics, counters etc... (which is implemented in the base class)

```
def reset(self) -> torch.Tensor:
    self._reset()
    self._df = self._df.sample(frac=1).reset_index(drop=True)
    return self._get_context()
```

The new bandit does not explicitly need to implement the `step` method since this is already implemnted in the base class. We do however need to implement `_compute_reward` and `_get_context` which `step` uses.

In `_compute_reward`, we need to figure out whether the given action corresponds to the correct label for this index or not and return the reward appropriately. This method also return the maximum possible reward in the current context which is used to compute regret.

```
def _compute_reward(self, action: int) -> Tuple[int, int]:
    label = self._df.iloc[self.idx, 0]
    r = int(label == (action + 1))
    return r, 1
```

The `_get_context` method should return a 13-dimensional `torch.Tensor` (in this case) corresponding to the context for the current index.

```
def _get_context(self) -> torch.Tensor:
    return torch.tensor(
        self._df.iloc[self.idx, 0].values,
        device=self.device,
        dtype=torch.float,
    )
```

Once you are done with the above, you can use the `WineDataBandit` class like you would any other bandit from `from genrl.utils.data_bandits`. You can use it with any of the `cb_agents` as well as training on it with `genrl.bandit.DCBTrainer`.

## Adding a new Deep Contextual Bandit Agent

The bandit submodule like all of `genrl` has been designed to be easily extensible for custom additions. This tutorial will show how to create a deep contextual bandit agent which will work with the rest of `genrl.bandit`

For the purpose of this tutorial we will consider a simple neural network based agent. Although this is a simplistic agent, implementation of any level of agent will need to have the following steps.

To start off with lets import necessary modules and make a class which inherits from `genrl.agents.bandits.contextual.base.DCBAgent`

```
from typing import Optional

import torch

from genrl.agents.bandits.contextual.base import DCBAgent
from genrl.agents.bandits.contextual.common import NeuralBanditModel, TransitionDB
from genrl.utils.data_bandits.base import DataBasedBandit

class NeuralAgent(DCBAgent):
    """Deep contextual bandit agent based on a neural network."""

    def __init__(self, bandit: DataBasedBandit, **kwargs):
        pass

    def select_action(self, context: torch.Tensor) -> int:
        pass

    def update_db(self, context: torch.Tensor, action: int, reward: int):
        pass

    def update_params(
            self,
            action: Optional[int] = None,
            batch_size: int = 512,
            train_epochs: int = 20,
    ):
        pass
```

We will need to implement `__init__`, `select_action`, `update_db` and `update_param` to make the class functional.

Lets start off with `__init__`. Here we will need to initialise some required parameters (`init_pulls`, `eval_with_dropout`, `t` and `update_count`) along with our transition database and the neural network. For the neural network, you can use the `NeuralBanditModel` class. It packages together many of the functionalities a neural network might require. Refer to the docs for more details.

```
def __init__(self, bandit: DataBasedBandit, **kwargs):
    super(NeuralAgent, self).__init__(bandit, **kwargs)
    self.model = (
        NeuralBanditModel(
            context_dim=self.context_dim,
            n_actions=self.n_actions,
            **kwargs
        )
        .to(torch.float)
        .to(self.device)
    )
    self.eval_with_dropout = kwargs.get("eval_with_dropout", False)
    self.db = TransitionDB(self.device)
    self.t = 0
    self.update_count = 0
```

For the select action function, the agent will pass the context vector through the neural network to produce logits for each action. It will then select the action with highest logit value. Note that it must also increment the timestep, and if take every action atleast `init_pulls` number of times initially.

```
def select_action(self, context: torch.Tensor) -> int:
    """Selects action for a given context"""
    self.model.use_dropout = self.eval_with_dropout
    self.t += 1
    if self.t < self.n_actions * self.init_pulls:
        return torch.tensor(
            self.t % self.n_actions, device=self.device, dtype=torch.int
        )

    results = self.model(context)
    action = torch.argmax(results["pred_rewards"]).to(torch.int)
    return action
```

For updating the databse we can use the `add` method of `TransitionDB` class.

```
def update_db(self, context: torch.Tensor, action: int, reward: int):
    """Updates transition database."""
    self.db.add(context, action, reward)
```

In `update_params` we need to train the model on the observations seen so far. Since the `NeuralBanditModel` class already hass a `train` function, we just need to call that. However if you are writing your own model, this is where the updates to the parameters would happen.

```
def update_params(
    self,
    action: Optional[int] = None,
    batch_size: int = 512,
    train_epochs: int = 20,
):
    """Update parameters of the agent."""
    self.update_count += 1
    self.model.train_model(self.db, train_epochs, batch_size)
```

Note that some of these functions have unused arguments. The signatures have been decided so as such to ensure generality over all classes of algorithms.

Once you are done with the above, you can use the `NeuralAgent` class like you would any other agent from `genrl.bandit`. You can use it with any of the bandits as well as training it with `genrl.bandit.DCBTrainer`.

## 2.3.2 Classical

### Q-Learning using GenRL

#### What is Q-Learning?

Q-Learning is one of the stepping stones for many reinforcement learning algorithms like DQN. AlphaGO is also one of the famous examples that use Q-Learning at the heart.

Essentially, a RL agent takes an action on the environment and then collects rewards and updates its policy, and over time gets better at collecting higher rewards.

In Q-Learning, we generally maintain a “Q-table” of *Q-values* by mapping them to a (state, action) pair.

A natural question is, What are these *Q-values*? It is nothing but the “Quality” of an action taken from a particular state. The more the *Q-value* the more chances of getting a better reward.

Q-Table is often initialized with random values/with zeros and as the agent collects rewards via performing actions on the environment we update this Q-Table at the  $i$  th step using the following formulation -

$$Q_i(s, a) = (1 - \alpha)Q_{i-1}(s, a) + \alpha * (\text{reward} + \gamma * \max_{a'} Q_{i-1}(s', a'))$$

Here  $\alpha$  is the learning rate in ML terms,  $\gamma$  is the discount factor for the rewards and  $s'$  is the state reached after taking action  $a$  from state  $s$ .

### FrozenLake-v0 environment

So to demonstrate how easy it is to train a Q-Learning approach in GenRL, we are taking a very simple gym environment.

Description of the environment (from the documentation) -

“The agent controls the movement of a character in a grid world. Some tiles of the grid are walkable, and others lead to the agent falling into the water. Additionally, the movement direction of the agent is uncertain and only partially depends on the chosen direction. The agent is rewarded for finding a walkable path to a goal tile.

Winter is here. You and your friends were tossing around a frisbee at the park when you made a wild throw that left the frisbee out in the middle of the lake. The water is mostly frozen, but there are a few holes where the ice has melted. If you step into one of those holes, you’ll fall into the freezing water. At this time, there’s an international frisbee shortage, so it’s absolutely imperative that you navigate across the lake and retrieve the disc. However, the ice is slippery, so you won’t always move in the direction you intend.

The surface is described using a grid like the following:

SFFF	(S: starting point, safe)
FHFH	(F: frozen surface, safe)
FFFH	(H: hole, fall to your doom)
HFFG	(G: goal, where the frisbee is located)

The episode ends when you reach the goal or fall in a hole. You receive a reward of 1 if you reach the goal, and zero otherwise.”

## Code

Let's import all the usefull stuff first.

```
import gym
from genrl import QLearning
from genrl.classical.common import Trainer
# for the agent
# for training the agent
```

Now that we have imported all the necessary stuff let's go ahead and define the environment, the agent and an object for the Trainer class.

```
env = gym.make("FrozenLake-v0")
agent = QLearning(env, gamma=0.6, lr=0.1, epsilon=0.1)
trainer = Trainer(
    agent,
    env,
    model="tabular",
    n_episodes=3000,
    start_steps=100,
    evaluate_frequency=100,
)
```

Great so far so good! Now moving towards the training process it is just calling the train method in the trainer class.

```
trainer.train()
trainer.evaluate()
```

That's it! You have successfully trained a Q-Learning agent. You can now go ahead and play with your own environments using GenRL!

## SARSA using GenRL

### What is SARSA?

SARSA is an acronym for State-Action-Reward-State-Action. It is an on-policy TD control method. Our aim is basically to estimate the Q-value or the utility value for state-action pair using the TD update rule given below.

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha * [R_{t+1} + \gamma * Q(S_{t+1}, A_{t+2}) - Q(S_t, A_t)]$$

### FrozenLake-v0 environment

So to demonstrate how easy it is to train a SARSA approach in GenRL, we are taking a very simple gym environment.

Description of the environment (from the documentation) -

“The agent controls the movement of a character in a grid world. Some tiles of the grid are walkable, and others lead to the agent falling into the water. Additionally, the movement direction of the agent is uncertain and only partially depends on the chosen direction. The agent is rewarded for finding a walkable path to a goal tile.

Winter is here. You and your friends were tossing around a frisbee at the park when you made a wild throw that left the frisbee out in the middle of the lake. The water is mostly frozen, but there are a few holes where the ice has melted. If you step into one of those holes, you'll fall into the freezing water. At this time, there's an international frisbee shortage, so it's absolutely imperative that you navigate across the lake and retrieve the disc. However, the ice is slippery, so you won't always move in the direction you intend.

The surface is described using a grid like the following:

SFFF	(S: starting point, safe)
FHFH	(F: frozen surface, safe)
FFFH	(H: hole, fall to your doom)
HFFG	(G: goal, where the frisbee is located)

The episode ends when you reach the goal or fall in a hole. You receive a reward of 1 if you reach the goal, and zero otherwise.”

## Code

Let’s import all the usefull stuff first.

```
import gym
from genrl import SARSA                                     # for the agent
from genrl.classical.common import Trainer                  # for training the agent
```

Now that we have imported all the necessary stuff let’s go ahead and define the environment, the agent and an object for the Trainer class.

```
env = gym.make("FrozenLake-v0")
agent = SARSA(env, gamma=0.6, lr=0.1, epsilon=0.1)
trainer = Trainer(
    agent,
    env,
    model="tabular",
    n_episodes=3000,
    start_steps=100,
    evaluate_frequency=100,
)
```

Great so far so good! Now moving towards the training process it is just calling the train method in the trainer class.

```
trainer.train()
trainer.evaluate()
```

That’s it! You have successfully trained a SARSA agent. You can now go ahead and play with your own environments using GenRL!

### 2.3.3 Deep RL Tutorials

#### Deep Reinforcement Learning Background

##### Background

The goal of Reinforcement Learning Algorithms is to maximize reward. This is usually achieved by having a policy  $\pi_\theta$  perform optimal behavior. Let’s denote this optimal policy by  $\pi_\theta^*$ . For ease, we define the Reinforcement Learning problem as a Markov Decision Process.

#### Markov Decision Process

An Markov Decision Process (MDP) is defined by  $(S, A, r, P_a)$  where,

- $S$  is a set of States.

- $A$  is a set of Actions.
- $r : S \rightarrow \mathbb{R}$  is a reward function.
- $P_a(s, s')$  is the transition probability that action  $a$  in state  $s$  leads to state  $s'$ .

Often we define two functions, a policy function  $\pi_\theta(s, a)$  and  $V_{\pi_\theta}(s)$ .

## Policy Function

The policy is the agent's strategy, we our goal is to make it optimal. The optimal policy is usually denoted by  $\pi_\theta^*$ . There are usually 2 types of policies:

### Stochastic Policy

The Policy Function is a stochastic variable defining a probability distribution over actions given states i.e. likelihood of every action when an agent is in a particular state. Formally,

$$\pi : S \times A \rightarrow [0, 1]$$

$$a \sim \pi(a|s)$$

### Deterministic Policy

The Policy Function maps from States directly to Actions.

$$\pi : S \rightarrow A$$

$$a = \pi(s)$$

## Value Function

The Value Function is defined as the expected return obtained when we follow a policy  $\pi$  starting from state  $S$ . Usually there are two types of value functions defined State Value Function and a State Action Value Function.

### State Value Function

The State Value Function is defined as the expected return starting from only State  $s$ .

$$V^\pi(s) = E[R_t]$$

### State Action Value Function

The Action Value Function is defined as the expected return starting from a state  $s$  and a taking an action  $a$ .

$$Q^\pi(s, a) = E[R_t]$$

The Action Value Function is also known as the **Quality** Function as it would denote how good a particular action is for a state  $s$ .

## Approximators

Neural Networks are often used as approximators for Policy and Value Functions. In such a case, we say these are **parameterised** by  $\theta$ . For e.g.  $\pi_\theta$ .

## Objective

The objective is to choose/learn a policy that will maximize a cumulative function of rewards received at each step, typically the discounted reward over a potential infinite horizon. We formulate this cumulative function as

$$E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

where we choose an action according to our policy,  $a_t = \pi_\theta(s_t)$ .

## Vanilla Policy Gradient

For background on Deep RL, its core definitions and problem formulations refer to Deep RL Background

## Objective

The objective is to choose/learn a policy that will maximize a cumulative function of rewards received at each step, typically the discounted reward over a potential infinite horizon. We formulate this cumulative function as

$$E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

where we choose the action  $a_t = \pi_\theta(s_t)$ .

## Algorithm Details

### Collect Experience

To make our agent learn, we first need to collect some experience in an online fashion. For this we make use of the `collect_rollouts` method. This method is defined in the `OnPolicyAgent` Base Class.

For updation, we would need to compute advantages from this experience. So, we store our experience in a Rollout Buffer. Action Selection —————

Note: We sample a **stochastic action** from the distribution on the action space by providing `False` as an argument to `select_action`.

For practical purposes we would assume that we are working with a finite horizon MDP.

### Update Equations

Let  $\pi_\theta$  denote a policy with parameters  $\theta$ , and  $J(\pi_\theta)$  denote the expected finite-horizon undiscounted return of the policy.

At each update timestep, we get value and log probabilities:

Now, that we have the log probabilities we calculate the gradient of  $J(\pi_\theta)$  as:

$$\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right],$$

where  $\tau$  is the trajectory.

We then update the policy parameters via stochastic gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_{\theta_k})$$

The key idea underlying vanilla policy gradients is to push up the probabilities of actions that lead to higher return, and push down the probabilities of actions that lead to lower return, until you arrive at the optimal policy.

## Training through the API

```
import gym

from genrl import VPG
from genrl.deep.common import OnPolicyTrainer
from genrl.environments import VectorEnv

env = VectorEnv("CartPole-v0")
agent = VPG('mlp', env)
trainer = OnPolicyTrainer(agent, env, log_mode=['stdout'])
trainer.train()
```

timestep	Episode	loss	mean_reward
0	0	9.1853	22.3825
20480	10	24.5517	80.3137
40960	20	24.4992	117.7011
61440	30	22.578	121.543
81920	40	20.423	114.7339
102400	50	21.7225	128.4013
122880	60	21.0566	116.034
143360	70	21.628	115.0562
163840	80	23.1384	133.4202
184320	90	23.2824	133.4202
204800	100	26.3477	147.87
225280	110	26.7198	139.7952
245760	120	30.0402	184.5045
266240	130	30.293	178.8646
286720	140	29.4063	162.5397
307200	150	30.9759	183.6771
327680	160	30.6517	186.1818
348160	170	31.7742	184.5045
368640	180	30.4608	186.1818
389120	190	30.2635	186.1818

## Advantage Actor Critic

For background on Deep RL, its core definitions and problem formulations refer to Deep RL Background

## Objective

The objective is to maximize the discounted cumulative reward function:

$$E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

This comprises of two parts in the Adantage Actor Critic Algorithm:

1. To choose/learn a policy that will increase the probability of landing an action that has higher expected return than the value of just the state and decrease the probability of landing an action that has lower expected return than the value of the state. The Advantage is computed as:

$$A(s, a) = Q(s, a) - V(s)$$

2. To learn a State Action Value Function (in the name of **Critic**) that estimates the future cumulative rewards given the current state and action. This function helps the policy in evaluation potential state, action pairs.

where we choose the action  $a_t = \pi_\theta(s_t)$ .

## Algorithm Details

### Action Selection and Values

ac here is an object of the `ActorCritic` class, which defined two methods: `get_value` and `get_action` and ofcourse they return the value approximation from the Critic and action from the Actor.

Note: We sample a **stochastic action** from the distribution on the action space by providing `False` as an argument to `select_action`.

For practical purposes we would assume that we are working with a finite horizon MDP.

### Collect Experience

To make our agent learn, we first need to collect some experience in an online fashion. For this we make use of the `collect_rollouts` method. This method is defined in the `OnPolicyAgent` Base Class.

For updation, we would need to compute advantages from this experience. So, we store our experience in a Rollout Buffer.

### Compute discounted Returns and Advantages

Next we can compute the advantages and the actual discounted returns for each state. This can be done very easily by simply calling `compute_returns_and_advantage`. Note this implementation of the rollout buffer is borrowed from Stable Baselines.

### Update Equations

Let  $\pi_\theta$  denote a policy with parameters  $\theta$ , and  $J(\pi_\theta)$  denote the expected finite-horizon undiscounted return of the policy.

At each update timestep, we get value and log probabilities:

Now, that we have the log probabilities we calculate the gradient of  $J(\pi_\theta)$  as:

$$\nabla_\theta J(\pi_\theta) = E_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) A^{\pi_\theta}(s_t, a_t) \right],$$

where  $\tau$  is the trajectory.

We then update the policy parameters via stochastic gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_{\theta_k})$$

The key idea underlying Advantage Actor Critic Algorithm is to push up the probabilities of actions that lead to higher return than the expected return of that state, and push down the probabilities of actions that lead to lower return than the expected return of that state, until you arrive at the optimal policy.

## Training through the API

```
import gym

from genrl import A2C
from genrl.deep.common import OnPolicyTrainer
from genrl.environments import VectorEnv

env = VectorEnv("CartPole-v0")
agent = A2C('mlp', env)
trainer = OnPolicyTrainer(agent, env, log_mode=['stdout'])
trainer.train()
```

## Proximal Policy Optimization

For background on Deep RL, its core definitions and problem formulations refer to Deep RL Background

## Objective

The objective is to maximize the discounted cumulative reward function:

$$E \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

The Proximal Policy Optimization Algorithm is very similar to the Advantage Actor Critic Algorithm except we add multiply the advantages with a ratio between the log probability of actions at experience collection time and at updaton time. What this does is - helps in establishing a trust region for not moving too away from the old policy and at the same time taking gradient ascent steps in the directions of actions which result in positive advantages.

where we choose the action  $a_t = \pi_\theta(s_t)$ .

## Algorithm Details

### Action Selection and Values

ac here is an object of the ActorCritic class, which defined two methods: `get_value` and `get_action` and ofcourse they return the value approximation from the Critic and action from the Actor.

Note: We sample a **stochastic action** from the distribution on the action space by providing `False` as an argument to `select_action`.

For practical purposes we would assume that we are working with a finite horizon MDP.

## Collect Experience

To make our agent learn, we first need to collect some experience in an online fashion. For this we make use of the `collect_rollouts` method. This method is defined in the `OnPolicyAgent` Base Class.

For updation, we would need to compute advantages from this experience. So, we store our experience in a Rollout Buffer.

## Compute discounted Returns and Advantages

Next we can compute the advantages and the actual discounted returns for each state. This can be done very easily by simply calling `compute_returns_and_advantage`. Note this implementation of the rollout buffer is borrowed from Stable Baselines.

## Update Equations

Let  $\pi_\theta$  denote a policy with parameters  $\theta$ , and  $J(\pi_\theta)$  denote the expected finite-horizon undiscounted return of the policy.

At each update timestep, we get value and log probabilities:

In the case of PPO our loss function is:

$$L(s, a, \theta_k, \theta) = \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \text{clip} \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_{\theta_k}}(s, a) \right),$$

where  $\tau$  is the trajectory.

We then update the policy parameters via stochastic gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_{\theta_k})$$

## Training through the API

```
import gym

from genrl import PPO1
from genrl.deep.common import OnPolicyTrainer
from genrl.environments import VectorEnv

env = VectorEnv("CartPole-v0")
agent = PPO1('mlp', env)
trainer = OnPolicyTrainer(agent, env, log_mode=['stdout'])
trainer.train()
```

### 2.3.4 Custom Policy Networks

GenRL provides default policies for images (CNNPolicy) and for other types of inputs(MlpPolicy). Sometimes, these default policies may be insufficient for your problem, or you may want more control over the policy definition, and hence require a custom policy.

The following code tutorial runs through the steps to use a custom policy depending on your problem.

Import the required libraries (eg. torch, torch.nn) and from GenRL, the algorithm (eg VPG), the trainer (eg. OnPolicyTrainer), the policy to be modified (eg. MlpPolicy)

```
# The necessary imports
import torch
import torch.nn as nn

from genrl import VPG
from genrl.core.policies import MlpPolicy
from genrl.environments import VectorEnv
from genrl.trainers import OnPolicyTrainer
```

Then define a `custom_policy` class that derives from the policy to be modified (in this case, the `MlpPolicy`)

```
# Define a custom MLP Policy
class custom_policy(MlpPolicy):
    def __init__(self, state_dim, action_dim, hidden, **kwargs):
        super().__init__(state_dim, action_dim, hidden)
        self.action_dim = action_dim
        self.state_dim = state_dim
```

The above class modifies the `MlpPolicy` to have the desired number of hidden layers in the MLP Neural network that parametrizes the policy. This is done by passing the variable `hidden` explicitly (`defaulthidden = (32, 32)`). The `state_dim` and `action_dim` variables stand for the dimensions of the state\_space and the action\_space, and are required to construct the neural network with the proper input and output shapes for your policy, given the environment.

In some cases, you may also want to redefine the policy used completely and not just customize an existing policy. This can be done by creating a new custom policy class that inherits the `BasePolicy` class. The `BasePolicy` class is a basic implementation of a general policy, with a `forward` and a `get_action` method. The `forward` method maps the input state to the action probabilities, and the `get_action` method selects an action from the given action probabilities (for both continuous and discrete action\_spaces)

Say you want to parametrize your policy by a Neural Network containing LSTM layers followed by MLP layers. This can be done as follows:

```
# Define a custom LSTM policy from the BasePolicy class
class custom_policy(BasePolicy):
    def __init__(self, state_dim, action_dim, hidden,
                 discrete=True, layer_size=512, layers=1, **kwargs):
        super(custom_policy, self).__init__(state_dim,
                                            action_dim,
                                            hidden,
                                            discrete,
                                            **kwargs)

        self.state_dim = state_dim
        self.action_dim = action_dim
        self.layer_size = layer_size
        self.lstm = nn.LSTM(self.state_dim, layer_size, layers)
        self.fc = mlp([layer_size] + list(hidden) + [action_dim],
                     sac=self.sac) # the mlp layers
```

(continues on next page)

(continued from previous page)

```
def forward(self, state):
    state, h = self.lstm(state.unsqueeze(0))
    state = state.view(-1, self.layer_size)
    action = self.fc(state)
    return action
```

Finally, it's time to train the custom policy. Define the environment to be trained on (CartPole-v0 in this case), and the state\_dim and action\_dim variables.

```
# Initialize an environment
env = VectorEnv("CartPole-v0", 1)

# Initialize the custom Policy
state_dim = env.observation_space.shape[0]
action_dim = env.action_space.n
policy = custom_policy(state_dim=state_dim, action_dim=action_dim,
                       hidden = (64, 64))
```

Then the algorithm is initialised with the custom policy defined, and the OnPolicyTrainer trains in with logging for better inference.

```
algo = VPG(policy, env)

# Initialize the trainer and start training
trainer = OnPolicyTrainer(algo, env, log_mode=["csv"],
                           logdir="./logs", epochs=100)
trainer.train()
```

## 2.3.5 Using A2C

### Using A2C on “CartPole-v0”

```
import gym

from genrl import A2C
from genrl.deep.common import OnPolicyTrainer
from genrl.environments import VectorEnv

env = VectorEnv("CartPole-v0")
agent = A2C('mlp', env, gamma=0.9, lr_policy=0.01, lr_value=0.1, policy_layers=(32,
                           32), value_layers=(32, 32), rollout_size=2048)
trainer = OnPolicyTrainer(agent, env, log_mode=['stdout', 'tensorboard'], log_key=
                           "Episode")
trainer.train()
```

### Using A2C on atari env - “Pong-v0”

```
env = VectorEnv("Pong-v0", env_type = "atari")
agent = A2C('cnn', env, gamma=0.99, lr_policy=0.01, lr_value=0.1, policy_layers=(32,
                           32), value_layers=(32, 32), rollout_size=2048)
```

(continues on next page)

(continued from previous page)

```
trainer = OnPolicyTrainer(agent, env, log_mode=['stdout', 'tensorboard'], log_key=
    ↪"timestep")
trainer.train()
```

More details can be found in the docs for [A2C](#) and [OnPolicyTrainer](#).

### 2.3.6 Vanilla Policy Gradient (VPG)

If you wanted to explore Policy Gradient algorithms in RL, there is a high chance you would've heard of PPO, DDPG, etc. but understanding them can be tricky if you're just starting.

VPG is arguably one of the easiest to understand policy gradient algorithms while still performing to a good enough level.

Let's understand policy gradient at a high level, unlike the classical algorithms like Q-Learning, Monte Carlo where you try to optimise the outputs of the action-value function of the agent which are then used to determine the optimal policy. In policy gradient, as one would like to say we go directly for the kill shot, basically we optimise the thing we want to use at the end, i.e. the Policy.

So that explains the “Policy” part of Policy Gradient, so what about “Gradient”, so gradient comes from the fact that we try to optimise the policy by gradient ascent (unlike the popular gradient descent, here we want to increase the values, hence ascent). So that explains the name, but how does it even work.

For that, have a look at the following Psuedo Code (source: [OpenAI](#))

For a more fundamental understanding [this](#) spinningup article is a good resource

Now that we have an understanding of how VPG works at a high level let's jump into the code to see it in action. This is a very minimal way to run a VPG agent on [GenRL](#)

#### VPG agent on a Cartpole Environment

```
import gym # OpenAI Gym

from genrl import VPG
from genrl.deep.common import OnPolicyTrainer
from genrl.environments import VectorEnv

env = VectorEnv("CartPole-v1")
agent = VPG('mlp', env)
trainer = OnPolicyTrainer(agent, env, epochs=200)
trainer.train()
```

This will run a VPG agent agent which will interact with the [CartPole-v1](#) gym environment Let's understand the output on running this (your individual values may differ),

timestep	Episode	loss	mean_reward
0	0	8.022	19.8835
20480	10	25.969	75.2941
40960	20	29.2478	144.2254
61440	30	25.5711	129.6203
81920	40	19.8718	96.6038
102400	50	19.2585	106.9452

(continues on next page)

(continued from previous page)

122880	60	17.7781	99.9024
143360	70	23.6839	121.543
163840	80	24.4362	129.2114
184320	90	28.1183	156.3359
204800	100	26.6074	155.1515
225280	110	27.2012	178.8646
245760	120	26.4612	164.498
266240	130	22.8618	148.4058
286720	140	23.465	153.4082
307200	150	21.9764	151.1439
327680	160	22.445	151.1439
348160	170	22.9925	155.7414
368640	180	22.6605	165.1613
389120	190	23.4676	177.316

timestep: It is basically the units of time the agent has interacted with the environment since the start of training  
 Episode: It is one complete rollout of the agent, to put it simply it is one complete run until the agent ends up winning or losing  
 loss: The loss encountered in that episode  
 mean\_reward: The mean reward accumulated in that episode

Now if you look closely the agent will not converge to the max reward even if you increase the epochs to say 5000, it is because that during training the agent is behaving according to a stochastic policy (Meaning when you try to pick from an action given a state from the policy it doesn't simply take the one with the maximum return, rather it samples an action from a probability distribution, so in other words, the policy isn't just like a lookup table, it's function which outputs a probability distribution over the actions which we sample from when using it to pick our optimal action). So even if the agent has figured out the optimal policy it is not taking the most optimal action at every step there is an inherent stochasticity to it. If we want the agent to make full use of the learnt policy we can add the following line of code at after the training

```
trainer.evaluate(render=True)
```

This will not only make the agent follow a deterministic policy and thus help you achieve the maximum reward possible reward attainable from the learnt policy but also allow you to see your agent perform by passing render=True

For more information on the VPG implementation and the various hyperparameters available have a look at the official GenRL docs [here](#)

Some more implementations

## VPG agent on an Atari Environment

```
env = VectorEnv("Pong-v0", env_type = "atari")
agent = VPG('cnn', env)
trainer = OnPolicyTrainer(agent, env, epochs=200)
trainer.train()
```

## 2.4 Agents

### 2.4.1 A2C

**genrl.agents.deep.a2c.a2c module**

```
class genrl.agents.deep.a2c.a2c.A2C(*args, noise: Any = None, noise_std: float = 0.1,
                                         value_coeff: float = 0.5, entropy_coeff: float = 0.01,
                                         **kwargs)
```

Bases: genrl.agents.deep.base.onpolicy.OnPolicyAgent

Advantage Actor Critic algorithm (A2C)

The synchronous version of A3C Paper: <https://arxiv.org/abs/1602.01783>

**network**

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

**env**

The environment that the agent is supposed to act on

**Type** Environment

**create\_model**

Whether the model of the algo should be created when initialised

**Type** bool

**batch\_size**

Mini batch size for loading experiences

**Type** int

**gamma**

The discount factor for rewards

**Type** float

**layers**

Layers in the Neural Network of the Q-value function

**Type** tuple of int

**lr\_policy**

Learning rate for the policy/actor

**Type** float

**lr\_value**

Learning rate for the critic

**Type** float

**rollout\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“rollout”]

**Type** str

**noise**

Action Noise function added to aid in exploration

**Type** ActionNoise

**noise\_std**  
Standard deviation of the action noise distribution  
**Type** float

**value\_coeff**  
Ratio of magnitude of value updates to policy updates  
**Type** float

**entropy\_coeff**  
Ratio of magnitude of entropy updates to policy updates  
**Type** float

**seed**  
Seed for randomness  
**Type** int

**render**  
Should the env be rendered during training?  
**Type** bool

**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]  
**Type** str

**empty\_logs()**  
Empties logs

**evaluate\_actions** (*states*: torch.Tensor, *actions*: torch.Tensor)  
Evaluates actions taken by actor  
Actions taken by actor and their respective states are analysed to get log probabilities and values from critics

**Parameters**

- **states** (torch.Tensor) – States encountered in rollout
- **actions** (torch.Tensor) – Actions taken in response to respective states

**Returns** Values of states encountered during the rollout log\_probs (torch.Tensor): Log of action probabilities given a state

**Return type** values (torch.Tensor)

**get\_hyperparams** () → Dict[str, Any]  
Get relevant hyperparameters to save  
**Returns** Hyperparameters to be saved  
**Return type** hyperparams (dict)

**get\_logging\_params** () → Dict[str, Any]  
Gets relevant parameters for logging  
**Returns** Logging parameters for monitoring training  
**Return type** logs (dict)

**get\_traj\_loss** (*values*: torch.Tensor, *dones*: torch.Tensor) → None  
Get loss from trajectory traversed by agent during rollouts

Computes the returns and advantages needed for calculating loss

#### Parameters

- **values** (`torch.Tensor`) – Values of states encountered during the rollout
- **dones** (`list of bool`) – Game over statuses of each environment

**load\_weights** (`weights`) → `None`

Load weights for the agent from pretrained model

#### Parameters **weights** (`dict`) – Dictionary of different neural net weights

**select\_action** (`state: numpy.ndarray, deterministic: bool = False`) → `numpy.ndarray`

Select action given state

Action Selection for On Policy Agents with Actor Critic

#### Parameters

- **state** (`np.ndarray`) – Current state of the environment
- **deterministic** (`bool`) – Should the policy be deterministic or stochastic

**Returns** Action taken by the agent value (`torch.Tensor`): Value of given state `log_prob` (`torch.Tensor`): Log probability of selected action

**Return type** `action (np.ndarray)`

**update\_params** () → `None`

Updates the the A2C network

Function to update the A2C actor-critic architecture

## 2.4.2 DDPG

### genrl.agents.deep.ddpg.ddpg module

**class** `genrl.agents.deep.ddpg.ddpg.DDPG(*args, noise: genrl.core.noise.ActionNoise = None, noise_std: float = 0.2, **kwargs)`  
Bases: `genrl.agents.deep.base.offpolicy.OffPolicyAgentAC`

Deep Deterministic Policy Gradient Algorithm

Paper: <https://arxiv.org/abs/1509.02971>

#### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

#### env

The environment that the agent is supposed to act on

**Type** Environment

#### create\_model

Whether the model of the algo should be created when initialised

**Type** bool

#### batch\_size

Mini batch size for loading experiences

**Type** int

**gamma**  
The discount factor for rewards  
**Type** float

**layers**  
Layers in the Neural Network of the Q-value function  
**Type** tuple of int

**lr\_policy**  
Learning rate for the policy/actor  
**Type** float

**lr\_value**  
Learning rate for the critic  
**Type** float

**replay\_size**  
Capacity of the Replay Buffer  
**Type** int

**buffer\_type**  
Choose the type of Buffer: [“push”, “prioritized”]  
**Type** str

**polyak**  
Target model update parameter (1 for hard update)  
**Type** float

**noise**  
Action Noise function added to aid in exploration  
**Type** ActionNoise

**noise\_std**  
Standard deviation of the action noise distribution  
**Type** float

**seed**  
Seed for randomness  
**Type** int

**render**  
Should the env be rendered during training?  
**Type** bool

**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]  
**Type** str

**empty\_logs()**  
Empties logs

**get\_hyperparams() → Dict[str, Any]**  
Get relevant hyperparameters to save  
**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

**get\_logging\_params ()** → Dict[str, Any]  
Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

**update\_params (update\_interval: int)** → None  
Update parameters of the model

**Parameters** **update\_interval** (int) – Interval between successive updates of the target model

### 2.4.3 DQN

#### genrl.agents.deep.dqn.base module

**class** genrl.agents.deep.dqn.base.DQN (\*args, max\_epsilon: float = 1.0, min\_epsilon: float = 0.01, epsilon\_decay: int = 1000, \*\*kwargs)  
Bases: genrl.agents.deep.base.offpolicy.OffPolicyAgent

Base DQN Class

Paper: <https://arxiv.org/abs/1312.5602>

##### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

##### env

The environment that the agent is supposed to act on

**Type** Environment

##### create\_model

Whether the model of the algo should be created when initialised

**Type** bool

##### batch\_size

Mini batch size for loading experiences

**Type** int

##### gamma

The discount factor for rewards

**Type** float

##### value\_layers

Layers in the Neural Network of the Q-value function

**Type** tuple of int

##### lr\_value

Learning rate for the Q-value function

**Type** float

##### replay\_size

Capacity of the Replay Buffer

**Type** int

**buffer\_type**  
Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**max\_epsilon**  
Maximum epsilon for exploration

**Type** str

**min\_epsilon**  
Minimum epsilon for exploration

**Type** str

**epsilon\_decay**  
Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**seed**  
Seed for randomness

**Type** int

**render**  
Should the env be rendered during training?

**Type** bool

**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**calculate\_epsilon\_by\_frame** () → float  
Helper function to calculate epsilon after every timestep  
Exponentially decays exploration rate from max epsilon to min epsilon. The greater the value of epsilon\_decay, the slower the decrease in epsilon

**empty\_logs** () → None  
Empties logs

**get\_greedy\_action** (state: torch.Tensor) → numpy.ndarray  
Greedy action selection

**Parameters** state (np.ndarray) – Current state of the environment

**Returns** Action taken by the agent

**Return type** action (np.ndarray)

**get\_hyperparams** () → Dict[str, Any]  
Get relevant hyperparameters to save

**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

**get\_logging\_params** () → Dict[str, Any]  
Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

**get\_q\_values** (*states*: torch.Tensor, *actions*: torch.Tensor) → torch.Tensor  
Get Q values corresponding to specific states and actions

**Parameters**

- **states** (torch.Tensor) – States for which Q-values need to be found
- **actions** (torch.Tensor) – Actions taken at respective states

**Returns** Q values for the given states and actions

**Return type** q\_values (torch.Tensor)

**get\_target\_q\_values** (*next\_states*: torch.Tensor, *rewards*: List[float], *dones*: List[bool]) → torch.Tensor  
Get target Q values for the DQN

**Parameters**

- **next\_states** (torch.Tensor) – Next states for which target Q-values need to be found
- **rewards** (list) – Rewards at each timestep for each environment
- **dones** (list) – Game over status for each environment

**Returns** Target Q values for the DQN

**Return type** target\_q\_values (torch.Tensor)

**load\_weights** (*weights*) → None  
Load weights for the agent from pretrained model

**Parameters** **weights** (Dict) – Dictionary of different neural net weights

**select\_action** (*state*: numpy.ndarray, *deterministic*: bool = False) → numpy.ndarray  
Select action given state

Epsilon-greedy action-selection

**Parameters**

- **state** (np.ndarray) – Current state of the environment
- **deterministic** (bool) – Should the policy be deterministic or stochastic

**Returns** Action taken by the agent

**Return type** action (np.ndarray)

**update\_params** (*update\_interval*: int) → None  
Update parameters of the model

**Parameters** **update\_interval** (int) – Interval between successive updates of the target model

**update\_params\_before\_select\_action** (*timestep*: int) → None  
Update necessary parameters before selecting an action

This updates the epsilon (exploration rate) of the agent every timestep

**Parameters** **timestep** (int) – Timestep of training

**update\_target\_model** () → None  
Function to update the target Q model

Updates the target model with the training model's weights when called

**genrl.agents.deep.dqn.categorical module**

```
class genrl.agents.deep.dqn.categorical.CategoricalDQN(*args, noisy_layers: Tuple =  
                                                    (32, 128), num_atoms: int =  
                                                    51, v_min: int = -10, v_max:  
                                                    int = 10, **kwargs)
```

Bases: *genrl.agents.deep.dqn.base.DQN*

Categorical DQN Algorithm

Paper: <https://arxiv.org/pdf/1707.06887.pdf>

**network**

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

**env**

The environment that the agent is supposed to act on

**Type** Environment

**create\_model**

Whether the model of the algo should be created when initialised

**Type** bool

**batch\_size**

Mini batch size for loading experiences

**Type** int

**gamma**

The discount factor for rewards

**Type** float

**layers**

Layers in the Neural Network of the Q-value function

**Type** tuple of int

**lr\_value**

Learning rate for the Q-value function

**Type** float

**replay\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**max\_epsilon**

Maximum epsilon for exploration

**Type** str

**min\_epsilon**

Minimum epsilon for exploration

**Type** str

**epsilon\_decay**

Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**noisy\_layers**

Noisy layers in the Neural Network of the Q-value function

**Type** tuple of int

**num\_atoms**

Number of atoms used in the discrete distribution

**Type** int

**v\_min**

Lower bound of value distribution

**Type** int

**v\_max**

Upper bound of value distribution

**Type** int

**seed**

Seed for randomness

**Type** int

**render**

Should the env be rendered during training?

**Type** bool

**device**

Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**get\_greedy\_action** (*state: torch.Tensor*) → numpy.ndarray

Greedy action selection

**Parameters** **state** (np.ndarray) – Current state of the environment

**Returns** Action taken by the agent

**Return type** action (np.ndarray)

**get\_q\_loss** (*batch: collections.namedtuple*)

Categorical DQN loss function to calculate the loss of the Q-function

**Parameters** **batch** (collections.namedtuple of torch.Tensor) – Batch of experiences

**Returns** Calculated loss of the Q-function

**Return type** loss (torch.Tensor)

**get\_q\_values** (*states: torch.Tensor, actions: torch.Tensor*)

Get Q values corresponding to specific states and actions

**Parameters**

- **states** (torch.Tensor) – States for which Q-values need to be found
- **actions** (torch.Tensor) – Actions taken at respective states

**Returns** Q values for the given states and actions

**Return type** `q_values (torch.Tensor)`

**get\_target\_q\_values** (`next_states: numpy.ndarray, rewards: List[float], dones: List[bool]`)

Projected Distribution of Q-values

Helper function for Categorical/Distributional DQN

#### Parameters

- **next\_states** (`torch.Tensor`) – Next states being encountered by the agent
- **rewards** (`torch.Tensor`) – Rewards received by the agent
- **dones** (`torch.Tensor`) – Game over status of each environment

**Returns** Projected Q-value Distribution or Target Q Values

**Return type** `target_q_values (object)`

## genrl.agents.deep.dqn.double module

**class** `genrl.agents.deep.dqn.double.DoubleDQN(*args, **kwargs)`

Bases: `genrl.agents.deep.dqn.base.DQN`

Double DQN Class

Paper: <https://arxiv.org/abs/1509.06461>

#### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

#### env

The environment that the agent is supposed to act on

**Type** Environment

#### batch\_size

Mini batch size for loading experiences

**Type** int

#### gamma

The discount factor for rewards

**Type** float

#### layers

Layers in the Neural Network of the Q-value function

**Type** tuple of int

#### lr\_value

Learning rate for the Q-value function

**Type** float

#### replay\_size

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**max\_epsilon**

Maximum epsilon for exploration

**Type** str

**min\_epsilon**

Minimum epsilon for exploration

**Type** str

**epsilon\_decay**

Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**seed**

Seed for randomness

**Type** int

**render**

Should the env be rendered during training?

**Type** bool

**device**

Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**get\_target\_q\_values** (*next\_states*: torch.Tensor, *rewards*: torch.Tensor, *dones*: torch.Tensor) →

torch.Tensor

Get target Q values for the DQN

**Parameters**

- **next\_states** (torch.Tensor) – Next states for which target Q-values need to be found
- **rewards** (list) – Rewards at each timestep for each environment
- **dones** (list) – Game over status for each environment

**Returns** Target Q values for the DQN

**Return type** target\_q\_values (torch.Tensor)

## genrl.agents.deep.dqn.dueling module

**class** genrl.agents.deep.dqn.dueling.**DuelingDQN** (\*args, \*\*kwargs)

Bases: *genrl.agents.deep.dqn.base.DQN*

Dueling DQN class

Paper: <https://arxiv.org/abs/1511.06581>

**network**

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

**env**

The environment that the agent is supposed to act on

**Type** Environment

**batch\_size**

Mini batch size for loading experiences

**Type** int

**gamma**

The discount factor for rewards

**Type** float

**layers**

Layers in the Neural Network of the Q-value function

**Type** tuple of int

**lr\_value**

Learning rate for the Q-value function

**Type** float

**replay\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**max\_epsilon**

Maximum epsilon for exploration

**Type** str

**min\_epsilon**

Minimum epsilon for exploration

**Type** str

**epsilon\_decay**

Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**seed**

Seed for randomness

**Type** int

**render**

Should the env be rendered during training?

**Type** bool

**device**

Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**genrl.agents.deep.dqn.noisy module**

**class** genrl.agents.deep.dqn.noisy.NoisyDQN(\*args, noisy\_layers: Tuple = (128, 128), \*\*kwargs)

Bases: *genrl.agents.deep.dqn.base.DQN*

Noisy DQN Algorithm

Paper: <https://arxiv.org/abs/1706.10295>

**network**

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

**env**

The environment that the agent is supposed to act on

**Type** Environment

**batch\_size**

Mini batch size for loading experiences

**Type** int

**gamma**

The discount factor for rewards

**Type** float

**layers**

Layers in the Neural Network of the Q-value function

**Type** tuple of int

**lr\_value**

Learning rate for the Q-value function

**Type** float

**replay\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**max\_epsilon**

Maximum epsilon for exploration

**Type** str

**min\_epsilon**

Minimum epsilon for exploration

**Type** str

**epsilon\_decay**

Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**noisy\_layers**

Noisy layers in the Neural Network of the Q-value function

**Type** tuple of int

**seed**  
Seed for randomness

**Type** int

**render**  
Should the env be rendered during training?

**Type** bool

**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

### genrl.agents.deep.dqn.prioritized module

**class** genrl.agents.deep.dqn.prioritized.**PrioritizedReplayDQN**(\*args, alpha: float = 0.6, beta: float = 0.4, \*\*kwargs)

Bases: *genrl.agents.deep.dqn.base.DQN*

Prioritized Replay DQN Class

Paper: <https://arxiv.org/abs/1511.05952>

**network**

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

**env**

The environment that the agent is supposed to act on

**Type** Environment

**batch\_size**

Mini batch size for loading experiences

**Type** int

**gamma**

The discount factor for rewards

**Type** float

**layers**

Layers in the Neural Network of the Q-value function

**Type** tuple of int

**lr\_value**

Learning rate for the Q-value function

**Type** float

**replay\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

---

**Type** str

**max\_epsilon**  
Maximum epsilon for exploration

**Type** str

**min\_epsilon**  
Minimum epsilon for exploration

**Type** str

**epsilon\_decay**  
Rate of decay of epsilon (in order to decrease exploration with time)

**Type** str

**alpha**  
Prioritization constant

**Type** float

**beta**  
Importance Sampling bias

**Type** float

**seed**  
Seed for randomness

**Type** int

**render**  
Should the env be rendered during training?

**Type** bool

**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**get\_q\_loss** (*batch: collections.namedtuple*) → torch.Tensor  
Normal Function to calculate the loss of the Q-function

**Parameters** **batch** (*collections.namedtuple* of *torch.Tensor*) – Batch of experiences

**Returns** Calculated loss of the Q-function

**Return type** loss (*torch.Tensor*)

## genrl.agents.deep.dqn.utils module

```
genrl.agents.deep.dqn.utils.categorical_greedy_action(agent:
                                                       genrl.agents.deep.dqn.base.DQN,
                                                       state: torch.Tensor) →
                                                       numpy.ndarray
```

Greedy action selection for Categorical DQN

### Parameters

- **agent** (*DQN*) – The agent
- **state** (*np.ndarray*) – Current state of the environment

**Returns** Action taken by the agent

**Return type** action (np.ndarray)

genrl.agents.deep.dqn.utils.**categorical\_q\_loss** (agent: genrl.agents.deep.dqn.base.DQN,  
batch: collections.namedtuple)

Categorical DQN loss function to calculate the loss of the Q-function

#### Parameters

- **agent** (DQN) – The agent

- **batch** (collections.namedtuple of torch.Tensor) – Batch of experiences

**Returns** Calculated loss of the Q-function

**Return type** loss (torch.Tensor)

genrl.agents.deep.dqn.utils.**categorical\_q\_target** (agent:  
genrl.agents.deep.dqn.base.DQN,  
next\_states: numpy.ndarray, rewards:  
List[float], dones: List[bool])

Projected Distribution of Q-values

Helper function for Categorical/Distributional DQN

#### Parameters

- **agent** (DQN) – The agent

- **next\_states** (torch.Tensor) – Next states being encountered by the agent

- **rewards** (torch.Tensor) – Rewards received by the agent

- **dones** (torch.Tensor) – Game over status of each environment

**Returns** Projected Q-value Distribution or Target Q Values

**Return type** target\_q\_values (object)

genrl.agents.deep.dqn.utils.**categorical\_q\_values** (agent:  
genrl.agents.deep.dqn.base.DQN,  
states: torch.Tensor, actions:  
torch.Tensor)

Get Q values given state for a Categorical DQN

#### Parameters

- **agent** (DQN) – The agent

- **states** (torch.Tensor) – States being replayed

- **actions** (torch.Tensor) – Actions being replayed

**Returns** Q values for the given states and actions

**Return type** q\_values (torch.Tensor)

genrl.agents.deep.dqn.utils.**ddqn\_q\_target** (agent: genrl.agents.deep.dqn.base.DQN,  
next\_states: torch.Tensor, rewards:  
torch.Tensor, dones: torch.Tensor) →  
torch.Tensor

Double Q-learning target

Can be used to replace the `get_target_values` method of the Base DQN class in any DQN algorithm

#### Parameters

- **agent** (DQN) – The agent
- **next\_states** (torch.Tensor) – Next states being encountered by the agent
- **rewards** (torch.Tensor) – Rewards received by the agent
- **dones** (torch.Tensor) – Game over status of each environment

**Returns** Target Q values using Double Q-learning

**Return type** target\_q\_values (torch.Tensor)

genrl.agents.deep.dqn.utils.**prioritized\_q\_loss** (*agent*: genrl.agents.deep.dqn.base.DQN,  
*batch*: collections.namedtuple)

Function to calculate the loss of the Q-function

**Returns** The agent loss (torch.Tensor): Calculateed loss of the Q-function

**Return type** agent (DQN)

## 2.4.4 PPO1

### genrl.agents.deep.ppo1.ppo1 module

**class** genrl.agents.deep.ppo1.ppo1.PPO1 (\*args, clip\_param: float = 0.2, value\_coeff: float = 0.5, entropy\_coeff: float = 0.01, \*\*kwargs)

Bases: genrl.agents.deep.base.onpolicy.OnPolicyAgent

Proximal Policy Optimization algorithm (Clipped policy).

Paper: <https://arxiv.org/abs/1707.06347>

#### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

#### env

The environment that the agent is supposed to act on

**Type** Environment

#### create\_model

Whether the model of the algo should be created when initialised

**Type** bool

#### batch\_size

Mini batch size for loading experiences

**Type** int

#### gamma

The discount factor for rewards

**Type** float

#### layers

Layers in the Neural Network of the Q-value function

**Type** tuple of int

#### lr\_policy

Learning rate for the policy/actor

**Type** float  
**lr\_value**  
Learning rate for the Q-value function  
**Type** float  
**rollout\_size**  
Capacity of the Rollout Buffer  
**Type** int  
**buffer\_type**  
Choose the type of Buffer: [“rollout”]  
**Type** str  
**clip\_param**  
Epsilon for clipping policy loss  
**Type** float  
**value\_coeff**  
Ratio of magnitude of value updates to policy updates  
**Type** float  
**entropy\_coeff**  
Ratio of magnitude of entropy updates to policy updates  
**Type** float  
**seed**  
Seed for randomness  
**Type** int  
**render**  
Should the env be rendered during training?  
**Type** bool  
**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]  
**Type** str  
**empty\_logs ()**  
Empties logs  
**evaluate\_actions (states: torch.Tensor, actions: torch.Tensor)**  
Evaluates actions taken by actor  
Actions taken by actor and their respective states are analysed to get log probabilities and values from critics  
**Parameters**

- **states** (torch.Tensor) – States encountered in rollout
- **actions** (torch.Tensor) – Actions taken in response to respective states

**Returns** Values of states encountered during the rollout log\_probs (torch.Tensor): Log of action probabilities given a state  
**Return type** values (torch.Tensor)

---

**get\_hyperparams()** → Dict[str, Any]  
Get relevant hyperparameters to save

**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

**get\_logging\_params()** → Dict[str, Any]  
Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

**get\_traj\_loss(values, dones)**  
Get loss from trajectory traversed by agent during rollouts  
Computes the returns and advantages needed for calculating loss

**Parameters**

- **values** (torch.Tensor) – Values of states encountered during the rollout
- **dones** (list of bool) – Game over statuses of each environment

**load\_weights(weights)** → None  
Load weights for the agent from pretrained model

**Parameters** **weights** (dict) – Dictionary of different neural net weights

**select\_action(state: numpy.ndarray, deterministic: bool = False)** → numpy.ndarray  
Select action given state  
Action Selection for On Policy Agents with Actor Critic

**Parameters**

- **state** (np.ndarray) – Current state of the environment
- **deterministic** (bool) – Should the policy be deterministic or stochastic

**Returns** Action taken by the agent value (torch.Tensor): Value of given state log\_prob (torch.Tensor): Log probability of selected action

**Return type** action (np.ndarray)

**update\_params()**  
Updates the the A2C network  
Function to update the A2C actor-critic architecture

## 2.4.5 VPG

### genrl.agents.deep.vpg.vpg module

```
class genrl.agents.deep.vpg.vpg.VPG(*args, **kwargs)
Bases: genrl.agents.deep.base.onpolicy.OnPolicyAgent
Vanilla Policy Gradient algorithm
Paper https://papers.nips.cc/paper/1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation.pdf
```

**network (str): The network type of the Q-value function.** Supported types: [“cnn”, “mlp”]

env (Environment): The environment that the agent is supposed to act on  
create\_model (bool): Whether the model of the algo should be created when initialised  
batch\_size (int): Mini batch size for loading experiences  
gamma (float): The discount factor for rewards  
layers (tuple of int): Layers in the Neural Network

of the Q-value function

lr\_policy (float): Learning rate for the policy/actor lr\_value (float): Learning rate for the Q-value function  
rollout\_size (int): Capacity of the Rollout Buffer  
buffer\_type (str): Choose the type of Buffer: [“rollout”]  
seed (int): Seed for randomness  
render (bool): Should the env be rendered during training?  
device (str): Hardware being used for training. Options:

[“cuda” -> GPU, “cpu” -> CPU]

### **empty\_logs ()**

Empties logs

### **get\_hyperparams () → Dict[str, Any]**

Get relevant hyperparameters to save

**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

### **get\_log\_probs (states: torch.Tensor, actions: torch.Tensor)**

Get log probabilities of action values

Actions taken by actor and their respective states are analysed to get log probabilities

#### **Parameters**

- **states** (torch.Tensor) – States encountered in rollout
- **actions** (torch.Tensor) – Actions taken in response to respective states

**Returns** Log of action probabilities given a state

**Return type** log\_probs (torch.Tensor)

### **get\_logging\_params () → Dict[str, Any]**

Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

### **get\_traj\_loss (values, dones)**

Get loss from trajectory traversed by agent during rollouts

Computes the returns and advantages needed for calculating loss

#### **Parameters**

- **values** (torch.Tensor) – Values of states encountered during the rollout
- **dones** (list of bool) – Game over statuses of each environment

### **load\_weights (weights) → None**

Load weights for the agent from pretrained model

**Parameters** **weights** (dict) – Dictionary of different neural net weights

### **select\_action (state: numpy.ndarray, deterministic: bool = False) → numpy.ndarray**

Select action given state

## Action Selection for Vanilla Policy Gradient

### Parameters

- **state** (np.ndarray) – Current state of the environment
- **deterministic** (bool) – Should the policy be deterministic or stochastic

### Returns

Action taken by the agent value (torch.Tensor): Value of given state. In VPG, there is no critic

to find the value so we set this to a default 0 for convenience

log\_prob (torch.Tensor): Log probability of selected action

**Return type** action (np.ndarray)

**update\_params()** → None

Updates the the A2C network

Function to update the A2C actor-critic architecture

## 2.4.6 TD3

### genrl.agents.deep.td3.td3 module

```
class genrl.agents.deep.td3.td3.TD3(*args, policy_frequency: int = 2, noise: genrl.core.noise.ActionNoise = None, noise_std: float = 0.2, **kwargs)
```

Bases: genrl.agents.deep.base.offpolicy.OffPolicyAgentAC

Twin Delayed DDPG Algorithm

Paper: <https://arxiv.org/abs/1509.02971>

#### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

#### env

The environment that the agent is supposed to act on

**Type** Environment

#### create\_model

Whether the model of the algo should be created when initialised

**Type** bool

#### batch\_size

Mini batch size for loading experiences

**Type** int

#### gamma

The discount factor for rewards

**Type** float

#### policy\_layers

Neural network layer dimensions for the policy

**Type** tuple of int  
**value\_layers**  
Neural network layer dimensions for the critics  
**Type** tuple of int  
**lr\_policy**  
Learning rate for the policy/actor  
**Type** float  
**lr\_value**  
Learning rate for the critic  
**Type** float  
**replay\_size**  
Capacity of the Replay Buffer  
**Type** int  
**buffer\_type**  
Choose the type of Buffer: [“push”, “prioritized”]  
**Type** str  
**polyak**  
Target model update parameter (1 for hard update)  
**Type** float  
**policy\_frequency**  
Frequency of policy updates in comparison to critic updates  
**Type** int  
**noise**  
Action Noise function added to aid in exploration  
**Type** ActionNoise  
**noise\_std**  
Standard deviation of the action noise distribution  
**Type** float  
**seed**  
Seed for randomness  
**Type** int  
**render**  
Should the env be rendered during training?  
**Type** bool  
**device**  
Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]  
**Type** str  
**empty\_logs ()**  
Empties logs  
**get\_hyperparams () → Dict[str, Any]**  
Get relevant hyperparameters to save

**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

**get\_logging\_params()** → Dict[str, Any]

Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

**update\_params** (update\_interval: int) → None

Update parameters of the model

**Parameters** update\_interval (int) – Interval between successive updates of the target model

## 2.4.7 SAC

### genrl.agents.deep.sac.sac module

**class** genrl.agents.deep.sac.sac.SAC (\*args, alpha: float = 0.01, polyak: float = 0.995, entropy\_tuning: bool = True, \*\*kwargs)

Bases: genrl.agents.deep.base.offpolicy.OffPolicyAgentAC

Soft Actor Critic algorithm (SAC)

Paper: <https://arxiv.org/abs/1812.05905>

#### network

The network type of the Q-value function. Supported types: [“cnn”, “mlp”]

**Type** str

#### env

The environment that the agent is supposed to act on

**Type** Environment

#### create\_model

Whether the model of the algo should be created when initialised

**Type** bool

#### batch\_size

Mini batch size for loading experiences

**Type** int

#### gamma

The discount factor for rewards

**Type** float

#### policy\_layers

Neural network layer dimensions for the policy

**Type** tuple of int

#### value\_layers

Neural network layer dimensions for the critics

**Type** tuple of int

**lr\_policy**

Learning rate for the policy/actor

**Type** float

**lr\_value**

Learning rate for the critic

**Type** float

**replay\_size**

Capacity of the Replay Buffer

**Type** int

**buffer\_type**

Choose the type of Buffer: [“push”, “prioritized”]

**Type** str

**alpha**

Entropy factor

**Type** str

**polyak**

Target model update parameter (1 for hard update)

**Type** float

**entropy\_tuning**

True if entropy tuning should be done, False otherwise

**Type** bool

**seed**

Seed for randomness

**Type** int

**render**

Should the env be rendered during training?

**Type** bool

**device**

Hardware being used for training. Options: [“cuda” -> GPU, “cpu” -> CPU]

**Type** str

**empty\_logs()**

Empties logs

**get\_alpha\_loss(log\_probs)**

Calculate Entropy Loss

**Parameters** **log\_probs** (float) – Log probs

**get\_hyperparams() → Dict[str, Any]**

Get relevant hyperparameters to save

**Returns** Hyperparameters to be saved

**Return type** hyperparams (dict)

**get\_logging\_params() → Dict[str, Any]**

Gets relevant parameters for logging

**Returns** Logging parameters for monitoring training

**Return type** logs (dict)

**get\_p\_loss** (states: torch.Tensor) → torch.Tensor  
Function to get the Policy loss

**Parameters** **states** (torch.Tensor) – States for which Q-values need to be found

**Returns** Calculated policy loss

**Return type** loss (torch.Tensor)

**get\_target\_q\_values** (next\_states: torch.Tensor, rewards: List[float], dones: List[bool]) → torch.Tensor  
Get target Q values for the SAC

**Parameters**

- **next\_states** (torch.Tensor) – Next states for which target Q-values need to be found
- **rewards** (list) – Rewards at each timestep for each environment
- **dones** (list) – Game over status for each environment

**Returns** Target Q values for the SAC

**Return type** target\_q\_values (torch.Tensor)

**select\_action** (state: numpy.ndarray, deterministic: bool = False) → numpy.ndarray  
Select action given state

Action Selection

**Parameters**

- **state** (np.ndarray) – Current state of the environment
- **deterministic** (bool) – Should the policy be deterministic or stochastic

**Returns** Action taken by the agent

**Return type** action (np.ndarray)

**update\_params** (update\_interval: int) → None  
Update parameters of the model

**Parameters** **update\_interval** (int) – Interval between successive updates of the target model

**update\_target\_model** () → None  
Function to update the target Q model

Updates the target model with the training model's weights when called

## 2.4.8 Q-Learning

### genrl.agents.classical.qlearning module

```
class genrl.agents.classical.qlearning.QLearning(env: gym.core.Env,
                                                epsilon: float = 0.9,
                                                gamma: float = 0.95,
                                                lr: float = 0.01)
```

Bases: object

Q-Learning Algorithm.

Paper- <https://link.springer.com/article/10.1007/BF00992698>

**env**

Environment with which agent interacts.

**Type** gym.Env

**epsilon**

exploration coefficient for epsilon-greedy exploration.

**Type** float, optional

**gamma**

discount factor.

**Type** float, optional

**lr**

learning rate for optimizer.

**Type** float, optional

**get\_action** (*state: numpy.ndarray, explore: bool = True*) → numpy.ndarray

Epsilon greedy selection of epsilon in the explore phase.

**Parameters**

- **state** (*np.ndarray*) – Environment state.
- **explore** (*bool, optional*) – True if exploration is required. False if not.

**Returns** action.

**Return type** np.ndarray

**get\_hyperparams** () → Dict[str, Any]

**update** (*transition: Tuple*) → None

Update the Q table.

**Parameters** **transition** (*Tuple*) – transition 4-tuple used to update Q-table. In the form  
(*state, action, reward, next\_state*)

## 2.4.9 SARSA

### genrl.agents.classical.sarsa.sarsa module

**class** genrl.agents.classical.sarsa.sarsa.**SARSA** (*env: gym.core.Env, epsilon: float = 0.9, lmbda: float = 0.9, gamma: float = 0.95, lr: float = 0.01*)

Bases: object

SARSA Algorithm.

Paper- <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.17.2539&rep=rep1&type=pdf>

**env**

Environment with which agent interacts.

**Type** gym.Env

**epsilon**

exploration coefficient for epsilon-greedy exploration.

**Type** float, optional

**gamma**  
discount factor.

**Type** float, optional

**lr**  
learning rate for optimizer.

**Type** float, optional

**get\_action**(*state*: numpy.ndarray, *explore*: bool = True) → numpy.ndarray  
Epsilon greedy selection of epsilon in the explore phase.

**Parameters**

- **state** (np.ndarray) – Environment state.
- **explore** (bool, optional) – True if exploration is required. False if not.

**Returns** action.

**Return type** np.ndarray

**update**(*transition*: Tuple) → None  
Update the Q table and e values

**Parameters** **transition** (Tuple) – transition 4-tuple used to update Q-table. In the form (state, action, reward, next\_state)

## 2.4.10 Contextual Bandit

### Base

```
class genrl.agents.bandits.contextual.base.DCBAgent(bandit: genrl.core.bandit.Bandit,
                                                    device: str = 'cpu', **kwargs)
```

Bases: genrl.core.bandit.BanditAgent

Base class for deep contextual bandit solving agents

#### Parameters

- **bandit** (gennav.deep.bandit.data\_bandits.DataBasedBandit) – The bandit to solve
- **device** (str) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

#### bandit

The bandit to solve

**Type** gennav.deep.bandit.data\_bandits.DataBasedBandit

#### device

Device to use for tensor operations.

**Type** torch.device

**select\_action**(*context*: torch.Tensor) → int

Select an action based on given context

**Parameters** **context** (torch.Tensor) – The context vector to select action for

---

**Note:** This method needs to be implemented in the specific agent.

---

**Returns** The action to take

**Return type** int

**update\_parameters** (action: *Optional[int] = None*, batch\_size: *Optional[int] = None*, train\_epochs: *Optional[int] = None*) → None  
Update parameters of the agent.

**Parameters**

- **action** (*Optional[int], optional*) – Action to update the parameters for. Defaults to None.
- **batch\_size** (*Optional[int], optional*) – Size of batch to update parameters with. Defaults to None.
- **train\_epochs** (*Optional[int], optional*) – Epochs to train neural network for. Defaults to None.

---

**Note:** This method needs to be implemented in the specific agent.

---

## Bootstrap Neural

**class** `genrl.agents.bandits.contextual.bootstrap_neural.BootstrapNeuralAgent` (*bandit: genrl.utils.data\_bandit, \*\*kwargs*)

Bases: `genrl.agents.bandits.contextual.base.DCBAgent`

Bootstraped ensemble agentfor deep contextual bandits.

**Parameters**

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int, optional*) – Number of times to select each action initially. Defaults to 3.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network. Defaults to [50, 50].
- **init\_lr** (*float, optional*) – Initial learning rate. Defaults to 0.1.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate. Defaults to 0.5.
- **lr\_reset** (*bool, optional*) – Whether to reset learning rate ever train interval. Defaults to True.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping. Defaults to 0.5.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **eval\_with\_dropout** (*bool, optional*) – Whether or not to use dropout at inference. Defaults to False.
- **n** (*int, optional*) – Number of models in ensemble. Defaults to 10.

- **add\_prob** (*float, optional*) – Probability of adding a transition to a database. Defaults to 0.95.
- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → int

Select an action based on given context.

Selects an action by computing a forward pass through a randomly selected network from the ensemble.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take.

**Return type** int

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

The transition is added to each database with a certain probability.

**Parameters**

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: Optional[int] = None, batch\_size: int = 512, train\_epochs: int = 20*)

Update parameters of the agent.

Trains each neural network in the ensemble.

**Parameters**

- **action** (*Optional[int], optional*) – Action to update the parameters for. Not applicable in this agent. Defaults to None.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*int, optional*) – Epochs to train neural network for. Defaults to 20

## Fixed

```
class genrl.agents.bandits.contextual.fixed.FixedAgent(bandit:
    genrl.utils.data_bandits.base.DataBasedBandit,
    p: List[float] = None, device: str = 'cpu')
```

Bases: *genrl.agents.bandits.contextual.base.DCBAgent*

**select\_action** (*context: torch.Tensor*) → int

Select an action based on fixed probabilities.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for. In this agent, context vector is not considered.

**Returns** The action to take.

**Return type** int

**update\_db** (\*args, \*\*kwargs)

**update\_params** (\*args, \*\*kwargs)

## Linear Posterior

**class** genrl.agents.bandits.contextual.linpos.**LinearPosteriorAgent** (*bandit*:  
genrl.utils.data\_bandits.base.DataBase  
\*\*kwargs)

Bases: genrl.agents.bandits.contextual.base.DCBAgent

Deep contextual bandit agent using bayesian regression for posterior inference.

### Parameters

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int, optional*) – Number of times to select each action initially. Defaults to 3.
- **lambda\_prior** (*float, optional*) – Guassian prior for linear model. Defaults to 0.25.
- **a0** (*float, optional*) – Inverse gamma prior for noise. Defaults to 6.0.
- **b0** (*float, optional*) – Inverse gamma prior for noise. Defaults to 6.0.
- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → int

Select an action based on given context.

Selecting action with highest predicted reward computed through betas sampled from posterior.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take.

**Return type** int

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

**Parameters**

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: int, batch\_size: int = 512, train\_epochs: Optional[int] = None*)

Update parameters of the agent.

Updated the posterior over beta though bayesian regression.

**Parameters**

- **action** (*int*) – Action to update the parameters for.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*Optional[int], optional*) – Epochs to train neural network for. Not applicable in this agent. Defaults to None

## Neural Greedy

```
class genrl.agents.bandits.contextual.neural_greedy.NeuralGreedyAgent (bandit:
    genrl.utils.data_bandits.base.L
    **kwargs)
```

Bases: *genrl.agents.bandits.contextual.base.DCBAgent*

Deep contextual bandit agent using epsilon greedy with a neural network.

### Parameters

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int, optional*) – Number of times to select each action initially. Defaults to 3.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network. Defaults to [50, 50].
- **init\_lr** (*float, optional*) – Initial learning rate. Defaults to 0.1.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate. Defaults to 0.5.
- **lr\_reset** (*bool, optional*) – Whether to reset learning rate ever train interval. Defaults to True.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping. Defaults to 0.5.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **eval\_with\_dropout** (*bool, optional*) – Whether or not to use dropout at inference. Defaults to False.
- **epsilon** (*float, optional*) – Probability of selecting a random action. Defaults to 0.0.
- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → *int*

Select an action based on given context.

Selects an action by computing a forward pass through network with an epsilon probability of selecting a random action.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take.

**Return type** *int*

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

### Parameters

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: Optional[int] = None, batch\_size: int = 512, train\_epochs: int = 20*)

Update parameters of the agent.

Trains neural network.

#### Parameters

- **action** (*Optional[int], optional*) – Action to update the parameters for. Not applicable in this agent. Defaults to None.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*int, optional*) – Epochs to train neural network for. Defaults to 20

## Neural Linear Posterior

```
class genrl.agents.bandits.contextual.neural_linpos.NeuralLinearPosteriorAgent(bandit:  
                                genrl.utils.data_...  
                                **kwargs)
```

Bases: *genrl.agents.bandits.contextual.base.DCBAgent*

Deep contextual bandit agent using bayesian regression on for posterior inference

A neural network is used to transform context vector to a latent representation on which bayesian regression is performed.

#### Parameters

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int, optional*) – Number of times to select each action initially. Defaults to 3.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network. Defaults to [50, 50].
- **init\_lr** (*float, optional*) – Initial learning rate. Defaults to 0.1.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate. Defaults to 0.5.
- **lr\_reset** (*bool, optional*) – Whether to reset learning rate ever train interval. Defaults to True.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping. Defaults to 0.5.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **eval\_with\_dropout** (*bool, optional*) – Whether or not to use dropout at inference. Defaults to False.
- **nn\_update\_ratio** (*int, optional*) – . Defaults to 2.
- **lambda\_prior** (*float, optional*) – Guassian prior for linear model. Defaults to 0.25.
- **a0** (*float, optional*) – Inverse gamma prior for noise. Defaults to 3.0.
- **b0** (*float, optional*) – Inverse gamma prior for noise. Defaults to 3.0.

- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → *int*

Select an action based on given context.

Selects an action by computing a forward pass through network to output a representation of the context on which bayesian linear regression is performed to select an action.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take.

**Return type** *int*

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

Updates latent context and predicted rewards seperately.

**Parameters**

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: int, batch\_size: int = 512, train\_epochs: int = 20*)

Update parameters of the agent.

Trains neural network and updates bayesian regression parameters.

**Parameters**

- **action** (*int*) – Action to update the parameters for.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*int, optional*) – Epochs to train neural network for. Defaults to 20

## Neural Noise Sampling

```
class genrl.agents.bandits.contextual.neural_noise_sampling.NeuralNoiseSamplingAgent(bandit:  
genrl.ut  
**kwargs)
```

Bases: *genrl.agents.bandits.contextual.base.DCBAgent*

Deep contextual bandit agent with noise sampling for neural network parameters.

**Parameters**

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int, optional*) – Number of times to select each action initially. Defaults to 3.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network. Defaults to [50, 50].
- **init\_lr** (*float, optional*) – Initial learning rate. Defaults to 0.1.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate. Defaults to 0.5.

- **lr\_reset** (*bool, optional*) – Whether to reset learning rate over train interval. Defaults to True.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping. Defaults to 0.5.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **eval\_with\_dropout** (*bool, optional*) – Whether or not to use dropout at inference. Defaults to False.
- **noise\_std\_dev** (*float, optional*) – Standard deviation of sampled noise. Defaults to 0.05.
- **eps** (*float, optional*) – Small constant for bounding KL divergence of noise. Defaults to 0.1.
- **noise\_update\_batch\_size** (*int, optional*) – Batch size for updating noise parameters. Defaults to 256.
- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → int

Select an action based on given context.

Selects an action by adding noise to neural network parameters and the computing forward with the context vector as input.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take

**Return type** int

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

**Parameters**

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: Optional[int] = None, batch\_size: int = 512, train\_epochs: int = 20*)

Update parameters of the agent.

Trains each neural network in the ensemble.

**Parameters**

- **action** (*Optional[int], optional*) – Action to update the parameters for. Not applicable in this agent. Defaults to None.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*int, optional*) – Epochs to train neural network for. Defaults to 20

## Variational

```
class genrl.agents.bandits.contextual.variational.VariationalAgent(bandit:  
    genrl.utils.data_bandits.base.DataBandit,  
    **kwargs)
```

Bases: *genrl.agents.bandits.contextual.base.DCBAgent*

Deep contextual bandit agent using variation inference.

### Parameters

- **bandit** (*DataBasedBandit*) – The bandit to solve
- **init\_pulls** (*int*, *optional*) – Number of times to select each action initially. Defaults to 3.
- **hidden\_dims** (*List[int]*, *optional*) – Dimensions of hidden layers of network. Defaults to [50, 50].
- **init\_lr** (*float*, *optional*) – Initial learning rate. Defaults to 0.1.
- **lr\_decay** (*float*, *optional*) – Decay rate for learning rate. Defaults to 0.5.
- **lr\_reset** (*bool*, *optional*) – Whether to reset learning rate ever train interval. Defaults to True.
- **max\_grad\_norm** (*float*, *optional*) – Maximum norm of gradients for gradient clipping. Defaults to 0.5.
- **dropout\_p** (*Optional[float]*, *optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **eval\_with\_dropout** (*bool*, *optional*) – Whether or not to use dropout at inference. Defaults to False.
- **noise\_std** (*float*, *optional*) – Standard deviation of noise in bayesian neural network. Defaults to 0.1.
- **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda. Defaults to “cpu”.

**select\_action** (*context: torch.Tensor*) → *int*

Select an action based on given context.

Selects an action by computing a forward pass through the bayesian neural network.

**Parameters** **context** (*torch.Tensor*) – The context vector to select action for.

**Returns** The action to take.

**Return type** *int*

**update\_db** (*context: torch.Tensor, action: int, reward: int*)

Updates transition database with given transition

### Parameters

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**update\_params** (*action: int, batch\_size: int = 512, train\_epochs: int = 20*)

Update parameters of the agent.

Trains each neural network in the ensemble.

### Parameters

- **action** (*Optional[int], optional*) – Action to update the parameters for. Not applicable in this agent. Defaults to None.
- **batch\_size** (*int, optional*) – Size of batch to update parameters with. Defaults to 512
- **train\_epochs** (*int, optional*) – Epochs to train neural network for. Defaults to 20

## 2.4.11 Multi-Armed Bandit

### Base

```
class genrl.agents.bandits.multiarmed.base.MABAgent(bandit:  
                                                    genrl.core.bandit.MultiArmedBandit)
```

Bases: genrl.core.bandit.BanditAgent

Base Class for Contextual Bandit solving Policy

#### Parameters

- **bandit** (*MultiArmedBandit type object*) – The Bandit to solve
- **requires\_init\_run** – Indicated if initialisation of Q values is required

#### action\_hist

Get the history of actions taken for contexts

**Returns** List of context, actions pairs

**Return type** list

#### counts

Get the number of times each action has been taken

**Returns** Numpy array with count for each action

**Return type** numpy.ndarray

#### regret

Get the current regret

**Returns** The current regret

**Return type** float

#### regret\_hist

Get the history of regrets incurred for each step

**Returns** List of rewards

**Return type** list

#### reward\_hist

Get the history of rewards received for each step

**Returns** List of rewards

**Return type** list

**select\_action** (*context: int*) → *int*

Select an action

This method needs to be implemented in the specific policy.

**Parameters** **context** (*int*) – the context to select action for

**Returns** Selected action

**Return type** *int*

**update\_params** (*context: int, action: int, reward: Union[int, float]*) → *None*

Update parameters for the policy

This method needs to be implemented in the specific policy.

**Parameters**

- **context** (*int*) – context for which action is taken
- **action** (*int*) – action taken for the step
- **reward** (*int or float*) – reward obtained for the step

## Bayesian Bandit

```
class genrl.agents.bandits.multiarmed.bayesian.BayesianUCBMAgent (bandit:  
    genrl.core.bandit.MultiArmedBandit)  
    alpha:  
        float =  
        1.0, beta:  
        float =  
        1.0, confidence:  
        float =  
        3.0)
```

Bases: *genrl.agents.bandits.multiarmed.base.MABAgent*

Multi-Armed Bandit Solver with Bayesian Upper Confidence Bound based Action Selection Strategy.

Refer to Section 2.7 of Reinforcement Learning: An Introduction.

**Parameters**

- **bandit** (*MultiArmedBandit type object*) – The Bandit to solve
- **alpha** (*float*) – alpha value for beta distribution
- **beta** (*float*) – beta values for beta distribution
- **c** (*float*) – Confidence level which controls degree of exploration

**a**

alpha parameter of beta distribution associated with the policy

**Type** numpy.ndarray

**b**

beta parameter of beta distribution associated with the policy

**Type** numpy.ndarray

**confidence**

Confidence level which weights the exploration term

**Type** float

**quality**

Q values for all the actions for alpha, beta and c

**Type** numpy.ndarray

**select\_action**(context: int) → int

Select an action according to bayesian upper confidence bound

Take action that maximises a weighted sum of the Q values and a beta distribution parameterized by alpha and beta and weighted by c for each action

**Parameters**

- **context** (int) – the context to select action for
- **t** (int) – timestep to choose action for

**Returns** Selected action

**Return type** int

**update\_params**(context: int, action: int, reward: float) → None

Update parameters for the policy

Updates the regret as the difference between max Q value and that of the action. Updates the Q values according to the reward received in this step

**Parameters**

- **context** (int) – context for which action is taken
- **action** (int) – action taken for the step
- **reward** (float) – reward obtained for the step

## Bernoulli Bandit

```
class genrl.agents.bandits.multiarmed.bernoulli_mab.BernoulliMAB(bandits: int
= 1, arms: int = 5, reward_probs: numpy.ndarray = None, context_type: str
= 'tensor')
```

Bases: genrl.core.bandit.MultiArmedBandit

Contextual Bandit with categorial context and bernoulli reward distribution

**Parameters**

- **bandits** (int) – Number of bandits
- **arms** (int) – Number of arms in each bandit
- **reward\_probs** (numpy.ndarray) – Probabilities of getting rewards

## Epsilon Greedy

```
class genrl.agents.bandits.multiarmed.epsgreedy.EpsGreedyMABAgent (bandit:  
                                     genrl.core.bandit.MultiArmedBandit  
                                     eps: float =  
                                     0.05)
```

Bases: *genrl.agents.bandits.multiarmed.base.MABAgent*

Contextual Bandit Policy with Epsilon Greedy Action Selection Strategy.

Refer to Section 2.3 of Reinforcement Learning: An Introduction.

### Parameters

- **bandit** (*MultiArmedBandit* type object) – The Bandit to solve
- **eps** (*float*) – Probability with which a random action is to be selected.

#### **eps**

Exploration constant

**Type** float

#### **quality**

Q values assigned by the policy to all actions

**Type** numpy.ndarray

#### **select\_action** (*context*: int) → int

Select an action according to epsilon greedy startegy

A random action is selected with epsilon probability over the optimal action according to the current Q values to encourage exploration of the policy.

**Parameters** **context** (int) – the context to select action for

**Returns** Selected action

**Return type** int

#### **update\_params** (*context*: int, *action*: int, *reward*: float) → None

Update parmeters for the policy

Updates the regret as the difference between max Q value and that of the action. Updates the Q values according to the reward recieived in this step.

### Parameters

- **context** (int) – context for which action is taken
- **action** (int) – action taken for the step
- **reward** (float) – reward obtained for the step

## Gaussian

```
class genrl.agents.bandits.multiarmed.gaussian_mab.GaussianMAB(bandits: int
= 10, arms: int = 5, reward_means: numpy.ndarray = None, context_type: str =
'tensor')
```

Bases: genrl.core.bandit.MultiArmedBandit

Contextual Bandit with categorial context and gaussian reward distribution

### Parameters

- **bandits** (*int*) – Number of bandits
- **arms** (*int*) – Number of arms in each bandit
- **reward\_means** (*numpy.ndarray*) – Mean of gaussian distribution for each reward

## Gradient

```
class genrl.agents.bandits.multiarmed.gradient.GradientMABAgent(bandit:
genrl.core.bandit.MultiArmedBandit,
alpha: float
= 0.1, temp:
float = 0.01)
```

Bases: *genrl.agents.bandits.multiarmed.base.MABAgent*

Multi-Armed Bandit Solver with Softmax Action Selection Strategy.

Refer to Section 2.8 of Reinforcement Learning: An Introduction.

### Parameters

- **bandit** (*MultiArmedBandit type object*) – The Bandit to solve
- **alpha** (*float*) – The step size parameter for gradient based update
- **temp** (*float*) – Temperature for softmax distribution over Q values of actions

#### alpha

Step size parameter for gradient based update of policy

**Type** float

#### probability\_hist

History of probability values assigned to each action for each timestep

**Type** numpy.ndarray

#### quality

Q values assigned by the policy to all actions

**Type** numpy.ndarray

#### select\_action(*context: int*) → int

Select an action according by softmax action selection strategy

Action is sampled from softmax distribution computed over the Q values for all actions

**Parameters** **context** (*int*) – the context to select action for

**Returns** Selected action

**Return type** int

**temp**

Temperature for softmax distribution over Q values of actions

**Type** float

**update\_params** (*context*: int, *action*: int, *reward*: float) → None

Update parameters for the policy

Updates the regret as the difference between max Q value and that of the action. Updates the Q values through a gradient ascent step

#### Parameters

- **context** (int) – context for which action is taken
- **action** (int) – action taken for the step
- **reward** (float) – reward obtained for the step

## Thompson Sampling

```
class genrl.agents.bandits.multiarmed.thompson.ThompsonSamplingMABAgent(bandit:
    genrl.core.bandit.MultiArmedBandit,
    alpha: float = 1.0,
    beta: float = 1.0)
```

Bases: *genrl.agents.bandits.multiarmed.base.MABAgent*

Multi-Armed Bandit Solver with Bayesian Upper Confidence Bound based Action Selection Strategy.

#### Parameters

- **bandit** (*MultiArmedBandit* type object) – The Bandit to solve
- **a** (float) – alpha value for beta distribution
- **b** (float) – beta values for beta distibution

**a**

alpha parameter of beta distribution associated with the policy

**Type** numpy.ndarray

**b**

beta parameter of beta distribution associated with the policy

**Type** numpy.ndarray

#### quality

Q values for all the actions for alpha, beta and c

**Type** numpy.ndarray

**select\_action** (*context: int*) → *int*

Select an action according to Thompson Sampling

Samples are taken from beta distribution parameterized by alpha and beta for each action. The action with the highest sample is selected.

**Parameters** **context** (*int*) – the context to select action for**Returns** Selected action**Return type** *int***update\_params** (*context: int, action: int, reward: float*) → *None*

Update parameters for the policy

Updates the regret as the difference between max Q value and that of the action. Updates the alpha value of beta distribution by adding the reward while the beta value is updated by adding 1 - reward. Update the counts the action taken.

**Parameters**

- **context** (*int*) – context for which action is taken
- **action** (*int*) – action taken for the step
- **reward** (*float*) – reward obtained for the step

## Upper Confidence Bound

**class** genrl.agents.bandits.multiarmed.ucb.**UCBMABAgent** (*bandit: genrl.core.bandit.MultiArmedBandit, confidence: float = 1.0*)Bases: *genrl.agents.bandits.multiarmed.base.MABAgent*

Multi-Armed Bandit Solver with Upper Confidence Bound based Action Selection Strategy.

Refer to Section 2.7 of Reinforcement Learning: An Introduction.

**Parameters**

- **bandit** (*MultiArmedBandit type object*) – The Bandit to solve
- **c** (*float*) – Confidence level which controls degree of exploration

**confidence**

Confidence level which weights the exploration term

**Type** *float***quality**

q values assigned by the policy to all actions

**Type** *numpy.ndarray***select\_action** (*context: int*) → *int*

Select an action according to upper confidence bound action selection

Take action that maximises a weighted sum of the Q values for the action and an exploration encouragement term controlled by c.

**Parameters** **context** (*int*) – the context to select action for**Returns** Selected action**Return type** *int*

**update\_params** (*context: int, action: int, reward: float*) → None

Update parameters for the policy

Updates the regret as the difference between max Q value and that of the action. Updates the Q values according to the reward received in this step.

#### Parameters

- **context** (*int*) – context for which action is taken
- **action** (*int*) – action taken for the step
- **reward** (*float*) – reward obtained for the step

## 2.5 Environments

### 2.5.1 Environments

#### Subpackages

#### Vectorized Environments

#### Submodules

#### genrl.environments.vec\_env.monitor module

```
class genrl.environments.vec_env.monitor.VecMonitor(venv:
                                                    genrl.environments.vec_env.vector_envs.VecEnv,
                                                    history_length: int = 0,
                                                    info_keys: Tuple = ())
Bases: genrl.environments.vec_env.wrappers.VecEnvWrapper
```

Monitor class for VecEnvs. Saves important variables into the info dictionary

#### Parameters

- **venv** (*object*) – Vectorized Environment
- **history\_length** (*int*) – Length of history for episode rewards and episode lengths
- **info\_keys** (*tuple or list*) – Important variables to save

**reset** () → numpy.ndarray

Resets Vectorized Environment

**Returns** Initial observations

**Return type** Numpy Array

**step** (*actions: numpy.ndarray*) → Tuple

Steps through all the environments and records important information

**Parameters** **actions** (*Numpy Array*) – Actions to be taken for the Vectorized Environment

**Returns** States, rewards, dones, infos

**genrl.environments.vec\_env.normalize module**

```
class genrl.environments.vec_env.normalize.VecNormalize(venv:  
                                                    genrl.environments.vec_env.vector_envs.VecEnv,  
                                                    norm_obs: bool = True,  
                                                    norm_reward: bool =  
                                                    True, clip_reward: float =  
                                                    20.0)  
Bases: genrl.environments.vec_env.wrappers.VecEnvWrapper
```

Wrapper to implement Normalization of observations and rewards for VecEnvs

**Parameters**

- **venv** (*Vectorized Environment*) – The Vectorized environment
- **n\_envs** (*int*) – Number of environments in VecEnv
- **norm\_obs** (*bool*) – True if observations should be normalized, else False
- **norm\_reward** (*bool*) – True if rewards should be normalized, else False
- **clip\_reward** (*float*) – Maximum absolute value for rewards

**close()**

Close all individual environments in the Vectorized Environment

**reset() → numpy.ndarray**

Resets Vectorized Environment

**Returns** Initial observations

**Return type** Numpy Array

**step(actions: numpy.ndarray) → Tuple**

Steps through all the environments and normalizes the observations and rewards (if enabled)

**Parameters** **actions** (*Numpy Array*) – Actions to be taken for the Vectorized Environment

**Returns** States, rewards, dones, infos

**genrl.environments.vec\_env.utils module**

```
class genrl.environments.vec_env.utils.RunningMeanStd(epsilon: float = 0.0001,  
                                                    shape: Tuple = ())
```

Bases: object

Utility Function to compute a running mean and variance calculator

**Parameters**

- **epsilon** (*float*) – Small number to prevent division by zero for calculations
- **shape** (*Tuple*) – Shape of the RMS object

**update(batch: numpy.ndarray)****genrl.environments.vec\_env.vector\_envs module**

```
class genrl.environments.vec_env.vector_envs.SerialVecEnv(*args, **kwargs)
```

Bases: genrl.environments.vec\_env.vector\_envs.VecEnv

Constructs a wrapper for serial execution through envs.

**close()**

Closes all envs

**get\_spaces()**

**images() → List[T]**

Returns an array of images from each env render

**render(mode='human')**

Renders all envs in a tiles format similar to baselines

**param mode** (Can either be ‘human’ or ‘rgb\_array’. Displays tiled

**images in ‘human’ and returns tiled images in ‘rgb\_array’)**

**type mode** string

**reset() → numpy.ndarray**

Resets all envs

**step(actions: numpy.ndarray) → Tuple**

Steps through all envs serially

**Parameters actions** (Iterable of ints/floats) – Actions from the model

**class genrl.environments.vec\_env.vector\_envs.SubProcessVecEnv(\*args, \*\*kwargs)**

Bases: [genrl.environments.vec\\_env.vector\\_envs.VecEnv](#)

Constructs a wrapper for parallel execution through envs.

**close()**

Closes all environments and processes

**get\_spaces() → Tuple**

Returns state and action spaces of environments

**reset() → numpy.ndarray**

Resets environments

**Returns** States after environment reset

**seed(seed: int = None)**

Sets seed for reproducability

**step(actions: numpy.ndarray) → Tuple**

Steps through environments serially

**Parameters actions** (Iterable of ints/floats) – Actions from the model

**class genrl.environments.vec\_env.vector\_envs.VecEnv(envs: List[T], n\_envs: int = 2)**

Bases: abc.ABC

Base class for multiple environments.

**Parameters**

- **env** (Gym Environment) – Gym environment to be vectorised
- **n\_envs** (int) – Number of environments

**action\_shape**

**action\_spaces**

```
close()
n_envs
obs_shape
observation_spaces
reset()
sample() → List[T]
    Return samples of actions from each environment
seed(seed: int)
    Set seed for reproducibility in all environments
step(actions)

genrl.environments.vec_env.vector_envs.Worker(parent_conn: multiprocessing.
                                                context.BaseContext.Pipe,
                                                child_conn: multiprocessing.
                                                context.BaseContext.Pipe, env:
                                                gym.core.Env)

Worker class to facilitate multiprocessing
```

#### Parameters

- **parent\_conn** (*Multiprocessing Pipe Connection*) – Parent connection of Pipe
- **child\_conn** (*Multiprocessing Pipe Connection*) – Child connection of Pipe
- **env** (*Gym Environment*) – Gym environment we need multiprocessing for

## genrl.environments.vec\_env.wrappers module

```
class genrl.environments.vec_env.wrappers.VecEnvWrapper(venv)
Bases: genrl.environments.vec_env.vector_envs.VecEnv

close()
render(mode='human')
reset()
step(actions)
```

## Module contents

### Submodules

## genrl.environments.action\_wrappers module

```
class genrl.environments.action_wrappers.ClipAction(env: Union[gym.core.Env,
                                                               genrl.environments.vec_env.vector_envs.VecEnv])
Bases: gym.core.ActionWrapper

Action Wrapper to clip actions

Parameters env (object) – The environment whose actions need to be clipped
action (action: numpy.ndarray) → numpy.ndarray
```

---

```
class genrl.environments.action_wrappers.RescaleAction(env: Union[gym.core.Env,
    genrl.environments.vec_env.vector_envs.VecEnv],
low: int, high: int)
```

Bases: gym.core.ActionWrapper

Action Wrapper to rescale actions

#### Parameters

- **env** (*object*) – The environment whose actions need to be rescaled
- **low** (*int*) – Lower limit of action
- **high** (*int*) – Upper limit of action

**action** (*action*: numpy.ndarray) → numpy.ndarray

## genrl.environments.atari\_preprocessing module

```
class genrl.environments.atari_preprocessing.AtariPreprocessing(env:
    gym.core.Env,
    frameskip:
    Union[Tuple,
    int] = (2, 5),
    grayscale:
    bool = True,
    screen_size:
    int = 84)
```

Bases: gym.core.Wrapper

Implementation for Image preprocessing for Gym Atari environments. Implements: 1) Frameskip 2) Grayscale 3) Downsampling to square image

**param env** Atari environment

**param frameskip** Number of steps between actions. E.g. frameskip=4 will mean 1 action will be taken for every 4 frames. It'll be a tuple

**if non-deterministic and a random number will be chosen from (2, 5)**

**param grayscale** Whether or not the output should be converted to grayscale

**param screen\_size** Size of the output screen (square output)

**type env** Gym Environment

**type frameskip** tuple or int

**type grayscale** boolean

**type screen\_size** int

**reset** () → numpy.ndarray

Resets state of environment

**Returns** Initial state

**Return type** NumPy array

**step** (*action*: numpy.ndarray) → numpy.ndarray

Step through Atari environment for given action

**Parameters** **action** (NumPy array) – Action taken by agent

**Returns** Current state, reward(for frameskip number of actions), done, info

### genrl.environments.atari\_wrappers module

**class** genrl.environments.atari\_wrappers.**FireReset** (*env: gym.core.Env*)  
Bases: gym.core.Wrapper

Some Atari environments do not actually do anything until a specific action (the fire action) is taken, so we make it take the action before starting the training process

**Parameters** **env** (*Gym Environment*) – Atari environment

**reset** () → numpy.ndarray

Resets state of environment. Performs the noop action a random number of times to introduce stochasticity

**Returns** Initial state

**Return type** NumPy array

**class** genrl.environments.atari\_wrappers.**NoopReset** (*env: gym.core.Env, max\_noops: int = 30*)  
Bases: gym.core.Wrapper

Some Atari environments always reset to the same state. So we take a random number of some empty (noop) action to introduce some stochasticity.

**Parameters**

- **env** (*Gym Environment*) – Atari environment
- **max\_noops** (*int*) – Maximum number of Noops to be taken

**reset** () → numpy.ndarray

Resets state of environment. Performs the noop action a random number of times to introduce stochasticity

**Returns** Initial state

**Return type** NumPy array

**step** (*action: numpy.ndarray*) → numpy.ndarray

Step through underlying Atari environment for given action

**Parameters** **action** (*NumPy array*) – Action taken by agent

**Returns** Current state, reward(for frameskip number of actions), done, info

### genrl.environments.base\_wrapper module

**class** genrl.environments.base\_wrapper.**BaseWrapper** (*env: Any, batch\_size: int = None*)  
Bases: abc.ABC

Base class for all wrappers

**batch\_size**

The number of batches trained per update

**close** () → None

Closes environment and performs any other cleanup

Must be overridden by subclasses

**render** () → None

Render the environment

---

**reset ()** → None  
Resets state of environment  
Must be overriden by subclasses

**Returns** Initial state

**seed (seed: int = None)** → None  
Set seed for environment

**step (action: numpy.ndarray)** → None  
Step through the environment  
Must be overriden by subclasses

## genrl.environments.frame\_stack module

**class** genrl.environments.frame\_stack.**FrameStack** (*env: gym.core.Env, framestack: int = 4, compress: bool = True*)

Bases: gym.core.Wrapper

Wrapper to stack the last few(4 by default) observations of agent efficiently

### Parameters

- **env** (*Gym Environment*) – Environment to be wrapped
- **framestack** (*int*) – Number of frames to be stacked
- **compress** (*bool*) – True if we want to use LZ4 compression to conserve memory usage

**reset ()** → numpy.ndarray

Resets environment

**Returns** Initial state of environment

**Return type** NumPy Array

**step (action: numpy.ndarray)** → numpy.ndarray

Steps through environment

**Parameters** **action** (*NumPy Array*) – Action taken by agent

**Returns** Next state, reward, done, info

**Return type** NumPy Array, float, boolean, dict

**class** genrl.environments.frame\_stack.**LazyFrames** (*frames: List[T], compress: bool = False*)

Bases: object

Efficient data structure to save each frame only once. Can use LZ4 compression to optimizer memory usage.

### Parameters

- **frames** (*collections.deque*) – List of frames that needs to converted to a LazyFrames data structure
- **compress** (*boolean*) – True if we want to use LZ4 compression to conserve memory usage

### shape

Returns dimensions of other object

**genrl.environments.gym\_wrapper module**

```
class genrl.environments.gym_wrapper.GymWrapper(env: gym.core.Env)
    Bases: gym.core.Wrapper

    Wrapper class for all Gym Environments

    Parameters
        • env (string) – Gym environment name
        • n_envs (None, int) – Number of environments. None if not vectorised
        • parallel (boolean) – If vectorised, should environments be run through serially or
            parallelly

    action_shape
    close() → None
        Closes environment

    obs_shape
    render(mode: str = 'human') → None
        Renders all envs in a tiles format similar to baselines.

        Parameters mode (string) – Can either be ‘human’ or ‘rgb_array’. Displays tiled images in
            ‘human’ and returns tiled images in ‘rgb_array’

    reset() → numpy.ndarray
        Resets environment

        Returns Initial state

    sample() → numpy.ndarray
        Shortcut method to directly sample from environment’s action space

        Returns Random action from action space

        Return type NumPy Array

    seed(seed: int = None) → None
        Set environment seed

        Parameters seed (int) – Value of seed

    step(action: numpy.ndarray) → numpy.ndarray
        Steps the env through given action

        Parameters action (NumPy array) – Action taken by agent

        Returns Next observation, reward, game status and debugging info
```

**genrl.environments.suite module**

```
genrl.environments.suite.AtariEnv(env_id: str, wrapper_list: List[T] = [<class
    'genrl.environments.atari_preprocessing.AtariPreprocessing',>
    <class 'genrl.environments.atari_wrappers.NoopReset',>
    <class 'genrl.environments.atari_wrappers.FireReset',>
    <class 'genrl.environments.time_limit.AtariTimeLimit',>
    <class 'genrl.environments.frame_stack.FrameStack'>])
    → gym.core.Env

Function to apply wrappers for all Atari envs by Trainer class
```

**Parameters**

- **env** (*string*) – Environment Name
- **wrapper\_list** (*list or tuple*) – List of wrappers to use

**Returns** Gym Atari Environment**Return type** object`genrl.environments.suite.GymEnv(env_id: str) → gym.core.Env`

Function to apply wrappers for all regular Gym envs by Trainer class

**Parameters** **env** (*string*) – Environment Name**Returns** Gym Environment**Return type** object`genrl.environments.suite.VectorEnv(env_id: str, n_envs: int = 2, parallel: int = False, env_type: str = 'gym') → genrl.environments.vec_env.vector_envs.VecEnv`

Chooses the kind of Vector Environment that is required

**param env\_id** Gym environment to be vectorised**param n\_envs** Number of environments**param parallel** True if we want environments to run parallelly and (**subprocesses, False if we want environments to run serially one after the other)****param env\_type** Type of environment. Currently, we support [“gym”, “atari”]**type env\_id** string**type n\_envs** int**type parallel** False**type env\_type** string**returns** Vector Environment**rtype** object**genrl.environments.time\_limit module**`class genrl.environments.time_limit.AtariTimeLimit(env, max_episode_len=None)`  
Bases: `gym.core.Wrapper`**reset** (\*\*kwargs)

Resets the state of the environment and returns an initial observation.

**Returns** the initial observation.**Return type** observation (object)**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)

```
class genrl.environments.time_limit.TimeLimit(env, max_episode_len=None)
```

Bases: gym.core.Wrapper

```
reset(**kwargs)
```

Resets the state of the environment and returns an initial observation.

**Returns** the initial observation.

**Return type** observation (object)

```
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)

## Module contents

## 2.6 Core

### 2.6.1 ActorCritic

```
class genrl.core.actor_critic.CNNActorCritic(framestack: int, action_dim: gym.spaces.space.Space, policy_layers: Tuple = (256, ), value_layers: Tuple = (256, ), val_type: str = 'V', discrete: bool = True, *args, **kwargs)
```

Bases: `genrl.core.base.BaseActorCritic`

CNN Actor Critic

**param framestack** Number of previous frames to stack together

**param action\_dim** Action dimensions of the environment

**param fc\_layers** Sizes of hidden layers

**param val\_type** Specifies type of value function: (

“V” for V(s), “Qs” for Q(s), “Qsa” for Q(s,a))

**param discrete** True if action space is discrete, else False

**param framestack** Number of previous frames to stack together

---

```

type action_dim int
type fc_layers tuple or list
type val_type str
type discrete bool

get_action (state: torch.Tensor, deterministic: bool = False) → torch.Tensor
    Get action from the Actor based on input
        param state The state being passed as input to the Actor
        param deterministic (True if the action space is deterministic,
            else False)
            type state Tensor
            type deterministic boolean
            returns action

get_value (inp: torch.Tensor) → torch.Tensor
    Get value from the Critic based on input
        Parameters inp (Tensor) – Input to the Critic
        Returns value

class genrl.core.actor_critic.MlpActorCritic (state_dim: gym.spaces.space.Space, action_dim: gym.spaces.space.Space, policy_layers: Tuple = (32, 32), value_layers: Tuple = (32, 32), val_type: str = 'V', discrete: bool = True, **kwargs)
    Bases: genrl.core.base.BaseActorCritic
    MLP Actor Critic

    state_dim
        State dimensions of the environment
            Type int

    action_dim
        Action space dimensions of the environment
            Type int

    hidden
        Hidden layers in the MLP
            Type list or tuple

    val_type
        Value type of the critic network
            Type str

    discrete
        True if the action space is discrete, else False
            Type bool

```

```
sac
    True if a SAC-like network is needed, else False
    Type bool

activation
    Activation function to be used. Can be either "tanh" or "relu"
    Type str

class genrl.core.actor_critic.MlpSingleActorMultiCritic(state_dim:
    gym.spaces.space.Space,
    action_dim:
    gym.spaces.space.Space,
    policy_layers: Tuple =
    (32, 32), value_layers:
    Tuple = (32, 32), val_type:
    str = 'V', discrete: bool =
    True, num_critics: int = 2,
    **kwargs)

Bases: genrl.core.base.BaseActorCritic
MLP Actor Critic

state_dim
    State dimensions of the environment
    Type int

action_dim
    Action space dimensions of the environment
    Type int

hidden
    Hidden layers in the MLP
    Type list or tuple

val_type
    Value type of the critic network
    Type str

discrete
    True if the action space is discrete, else False
    Type bool

num_critics
    Number of critics in the architecture
    Type int

sac
    True if a SAC-like network is needed, else False
    Type bool

activation
    Activation function to be used. Can be either "tanh" or "relu"
    Type str
```

---

**forward**(*x*)  
Defines the computation performed at every call.  
Should be overridden by all subclasses.

---

**Note:** Although the recipe for forward pass needs to be defined within this function, one should call the Module instance afterwards instead of this since the former takes care of running the registered hooks while the latter silently ignores them.

---

**get\_action**(*state*: *torch.Tensor*, *deterministic*: *bool* = *False*)  
Get action from the Actor based on input

**param state** The state being passed as input to the Actor  
**param deterministic** (True if the action space is deterministic,  
**else False**)  
**type state** Tensor  
**type deterministic** boolean  
**returns** action

**get\_value**(*state*: *torch.Tensor*, *mode*=’first’) → *torch.Tensor*  
Get Values from the Critic

**Arg:** *state* (*torch.Tensor*): The state(s) being passed to the critics mode (str): What values should be returned. Types:  
“both” → Both values will be returned “min” → The minimum of both values will be returned  
“first” → The value from the first critic only will be returned

**Returns** List of values as estimated by each individual critic  
**Return type** values (list)

genrl.core.actor\_critic.**get\_actor\_critic\_from\_name**(*name\_*: str)  
Returns Actor Critic given the type of the Actor Critic

**Parameters** **ac\_name** (str) – Name of the policy needed  
**Returns** Actor Critic class to be used

## 2.6.2 Base

**class** genrl.core.base.**BaseActorCritic**  
Bases: *torch.nn.modules.module.Module*  
Basic implementation of a general Actor Critic

**get\_action**(*state*: *torch.Tensor*, *deterministic*: *bool* = *False*) → *torch.Tensor*  
Get action from the Actor based on input

**param state** The state being passed as input to the Actor  
**param deterministic** (True if the action space is deterministic,  
**else False**)

```
    type state Tensor
    type deterministic boolean
    returns action

get_value(state: torch.Tensor) → torch.Tensor
    Get value from the Critic based on input

    Parameters state(Tensor) – Input to the Critic
    Returns value

class genrl.core.base.BasePolicy(state_dim: int, action_dim: int, hidden: Tuple, discrete: bool,
                                     **kwargs)
    Bases: torch.nn.modules.module.Module
    Basic implementation of a general Policy

    Parameters
        • state_dim(int) – State dimensions of the environment
        • action_dim(int) – Action dimensions of the environment
        • hidden(tuple or list) – Sizes of hidden layers
        • discrete(bool) – True if action space is discrete, else False

    forward(state: torch.Tensor) → Tuple[torch.Tensor, Optional[torch.Tensor]]
        Defines the computation performed at every call.

        Parameters state(Tensor) – The state being passed as input to the policy
    get_action(state: torch.Tensor, deterministic: bool = False) → torch.Tensor
        Get action from policy based on input
            param state The state being passed as input to the policy
            param deterministic (True if the action space is deterministic,
                else False)
                type state Tensor
                type deterministic boolean
                returns action

class genrl.core.base.BaseValue(state_dim: int, action_dim: int)
    Bases: torch.nn.modules.module.Module
    Basic implementation of a general Value function

    forward(state: torch.Tensor) → torch.Tensor
        Defines the computation performed at every call.

        Parameters state(Tensor) – Input to value function
    get_value(state: torch.Tensor) → torch.Tensor
        Get value from value function based on input
            Parameters state(Tensor) – Input to value function
            Returns Value
```

### 2.6.3 Buffers

```
class genrl.core.buffers.PrioritizedBuffer(capacity: int, alpha: float = 0.6, beta: float = 0.4)
```

Bases: object

Implements the Prioritized Experience Replay Mechanism

#### Parameters

- **capacity** (*int*) – Size of the replay buffer
- **alpha** (*int*) – Level of prioritization

**pos**

**push** (*inp: Tuple*) → None

Adds new experience to buffer

**param inp** (Tuple containing *state, action, reward,*

*next\_state and done*)

**type inp** tuple

**returns** None

**sample** (*batch\_size: int, beta: float = None*) → Tuple[torch.Tensor, torch.Tensor, torch.Tensor, torch.Tensor, torch.Tensor, torch.Tensor]

(Returns randomly sampled memories from replay memory along with their respective indices and weights)

**param batch\_size** Number of samples per batch

**param beta** (Bias exponent used to correct

**Importance Sampling (IS) weights**

**type batch\_size** int

**type beta** float

**returns** (Tuple containing *states, actions, next\_states,*

*rewards, dones, indices and weights*)

**update\_priorities** (*batch\_indices: Tuple, batch\_priorities: Tuple*) → None

Updates list of priorities with new order of priorities

**param batch\_indices** List of indices of batch

**param batch\_priorities** (List of priorities of the batch at the

specific indices)

**type batch\_indices** list or tuple

**type batch\_priorities** list or tuple

```
class genrl.core.buffers.PrioritizedReplayBufferSamples(states, actions, rewards,
                                                       next_states, dones, indices,
                                                       weights)
Bases: tuple

actions
    Alias for field number 1

dones
    Alias for field number 4

indices
    Alias for field number 5

next_states
    Alias for field number 3

rewards
    Alias for field number 2

states
    Alias for field number 0

weights
    Alias for field number 6

class genrl.core.buffers.PushReplayBuffer(capacity: int)
Bases: object

Implements the basic Experience Replay Mechanism

Parameters capacity (int) – Size of the replay buffer

push (inp: Tuple) → None
    Adds new experience to buffer

        Parameters inp (tuple) – Tuple containing state, action, reward, next_state and done

        Returns None

sample (batch_size: int) → Tuple[torch.Tensor, torch.Tensor, torch.Tensor, torch.Tensor, torch.Tensor]
    Returns randomly sampled experiences from replay memory

        param batch_size Number of samples per batch

        type batch_size int

        returns (Tuple composing of state, action, reward,
next_state and done)

class genrl.core.buffers.ReplayBuffer(size, env)
Bases: object

extend (inp)

push (inp)

sample (batch_size)

class genrl.core.buffers.ReplayBufferSamples(states, actions, rewards, next_states,
                                             dones)
Bases: tuple

actions
    Alias for field number 1
```

---

**dones**  
Alias for field number 4

**next\_states**  
Alias for field number 3

**rewards**  
Alias for field number 2

**states**  
Alias for field number 0

## 2.6.4 Noise

**class** genrl.core.noise.**ActionNoise** (*mean*: float, *std*: float)

Bases: abc.ABC

Base class for Action Noise

### Parameters

- **mean** (float) – Mean of noise distribution
- **std** (float) – Standard deviation of noise distribution

#### **mean**

Returns mean of noise distribution

#### **std**

Returns standard deviation of noise distribution

**class** genrl.core.noise.**NoisyLinear** (*in\_features*: int, *out\_features*: int, *std\_init*: float = 0.4)

Bases: torch.nn.modules.module.Module

Noisy Linear Layer Class

Class to represent a Noisy Linear class (noisy version of nn.Linear)

#### **in\_features**

Input dimensions

**Type** int

#### **out\_features**

Output dimensions

**Type** int

#### **std\_init**

Weight initialisation constant

**Type** float

#### **forward** (*state*: torch.Tensor) → torch.Tensor

Defines the computation performed at every call.

Should be overridden by all subclasses.

---

**Note:** Although the recipe for forward pass needs to be defined within this function, one should call the Module instance afterwards instead of this since the former takes care of running the registered hooks while the latter silently ignores them.

```
reset_noise() → None
    Reset noise components of layer

reset_parameters() → None
    Reset parameters of layer

class genrl.core.noise.NormalActionNoise(mean: float, std: float)
    Bases: genrl.core.noise.ActionNoise
        Normal implementation of Action Noise

    Parameters
        • mean (float) – Mean of noise distribution
        • std (float) – Standard deviation of noise distribution

    reset() → None

class genrl.core.noise.OrnsteinUhlenbeckActionNoise(mean: float, std: float, theta: float = 0.15, dt: float = 0.01, initial_noise: numpy.ndarray = None)
    Bases: genrl.core.noise.ActionNoise
        Ornstein Uhlenbeck implementation of Action Noise

    Parameters
        • mean (float) – Mean of noise distribution
        • std (float) – Standard deviation of noise distribution
        • theta (float) – Parameter used to solve the Ornstein Uhlenbeck process
        • dt (float) – Small parameter used to solve the Ornstein Uhlenbeck process
        • initial_noise (Numpy array) – Initial noise distribution

    reset() → None
        Reset the initial noise value for the noise distribution sampling
```

## 2.6.5 Policies

```
class genrl.core.policies.CNNPolicy(framestack: int, action_dim: int, hidden: Tuple = (32, 32), discrete: bool = True, *args, **kwargs)
    Bases: genrl.core.base.BasePolicy
        CNN Policy

    Parameters
        • framestack (int) – Number of previous frames to stack together
        • action_dim (int) – Action dimensions of the environment
        • fc_layers (tuple or list) – Sizes of hidden layers
        • discrete (bool) – True if action space is discrete, else False
        • channels (list or tuple) – Channel sizes for cnn layers

    forward(state: numpy.ndarray) → numpy.ndarray
        Defines the computation performed at every call.

    Parameters state (Tensor) – The state being passed as input to the policy
```

---

```
class genrl.core.policies.MlpPolicy(state_dim: int, action_dim: int, hidden: Tuple = (32, 32),  
                                     discrete: bool = True, *args, **kwargs)  
Bases: genrl.core.base.BasePolicy
```

MLP Policy

#### Parameters

- **state\_dim** (int) – State dimensions of the environment
- **action\_dim** (int) – Action dimensions of the environment
- **hidden** (tuple or list) – Sizes of hidden layers
- **discrete** (bool) – True if action space is discrete, else False

```
genrl.core.policies.get_policy_from_name(name_: str)
```

Returns policy given the name of the policy

**Parameters** **name** (str) – Name of the policy needed

**Returns** Policy Function to be used

## 2.6.6 RolloutStorage

```
class genrl.core.rollout_storage.BaseBuffer(buffer_size: int, env: Union[gym.core.Env,  
                                         genrl.environments.vec_env.vector_envs.VecEnv],  
                                         device: Union[torch.device, str] = 'cpu')
```

Bases: object

Base class that represent a buffer (rollout or replay) :param buffer\_size: (int) Max number of element in the buffer :param env: (Environment) The environment being trained on :param device: (Union[torch.device, str]) PyTorch device

to which the values will be converted

**Parameters** **n\_envs** – (int) Number of parallel environments

**add**(\*args, \*\*kwargs) → None  
Add elements to the buffer.

**extend**(\*args, \*\*kwargs) → None  
Add a new batch of transitions to the buffer

**reset**() → None  
Reset the buffer.

**sample**(*batch\_size*: int)

**Parameters** **batch\_size** – (int) Number of element to sample

**Returns** (Union[RolloutBufferSamples, ReplayBufferSamples])

**size**() → int

**Returns** (int) The current size of the buffer

**static swap\_and\_flatten**(*arr*: numpy.ndarray) → numpy.ndarray

Swap and then flatten axes 0 (*buffer\_size*) and 1 (*n\_envs*) to convert shape from [n\_steps, n\_envs, ...] (when ... is the shape of the features) to [n\_steps \* n\_envs, ...] (which maintain the order) :param arr: (np.ndarray) :return: (np.ndarray)

**to\_torch** (*array*: numpy.ndarray, *copy*: bool = True) → torch.Tensor  
Convert a numpy array to a PyTorch tensor. Note: it copies the data by default :param *array*: (np.ndarray)  
:param *copy*: (bool) Whether to copy or not the data  
(may be useful to avoid changing things be reference)

**Returns** (torch.Tensor)

**class** genrl.core.rollout\_storage.ReplayBufferSamples (*observations*, *actions*,  
*next\_observations*, *dones*,  
*rewards*)  
Bases: tuple

**actions**  
Alias for field number 1

**dones**  
Alias for field number 3

**next\_observations**  
Alias for field number 2

**observations**  
Alias for field number 0

**rewards**  
Alias for field number 4

**class** genrl.core.rollout\_storage.RolloutBuffer (*buffer\_size*: int, *env*:  
Union[gym.core.Env,  
genrl.environments.vec\_env.vector\_envs.VecEnv],  
*device*: Union[torch.device, str] = 'cpu',  
*gae\_lambda*: float = 1, *gamma*: float = 0.99)  
Bases: genrl.core.rollout\_storage.BaseBuffer

Rollout buffer used in on-policy algorithms like A2C/PPO. :param *buffer\_size*: (int) Max number of element in the buffer :param *env*: (Environment) The environment being trained on :param *device*: (torch.device) :param *gae\_lambda*: (float) Factor for trade-off of bias vs variance for Generalized Advantage Estimator

Equivalent to classic advantage when set to 1.

### Parameters

- **gamma** – (float) Discount factor
- **n\_envs** – (int) Number of parallel environments

**add** (*obs*: numpy.ndarray, *action*: numpy.ndarray, *reward*: numpy.ndarray, *done*: numpy.ndarray, *value*: torch.Tensor, *log\_prob*: torch.Tensor) → None

### Parameters

- **obs** – (np.ndarray) Observation
- **action** – (np.ndarray) Action
- **reward** – (np.ndarray)
- **done** – (np.ndarray) End of episode signal.
- **value** – (torch.Tensor) estimated value of the current state following the current policy.
- **log\_prob** – (torch.Tensor) log probability of the action following the current policy.

---

```

compute_returns_and_advantage (last_value: torch.Tensor, dones: numpy.ndarray, use_gae:  

    bool = False) → None  

Post-processing step: compute the returns (sum of discounted rewards) and advantage ( $A(s) = R - V(S)$ ).  

Adapted from Stable-Baselines PPO2. :param last_value: (torch.Tensor) :param dones: (np.ndarray)  

:param use_gae: (bool) Whether to use Generalized Advantage Estimation  

    or normal advantage for advantage computation.

get (batch_size: Optional[int] = None) → Generator[genrl.core.rollout_storage.RolloutBufferSamples,  

None, None]  

reset () → None  

    Reset the buffer.

class genrl.core.rollout_storage.RolloutBufferSamples (observations, actions,  

old_values, old_log_prob,  

advantages, returns)  

Bases: tuple

actions  

    Alias for field number 1

advantages  

    Alias for field number 4

observations  

    Alias for field number 0

old_log_prob  

    Alias for field number 3

old_values  

    Alias for field number 2

returns  

    Alias for field number 5

class genrl.core.rollout_storage.RolloutReturn (episode_reward, episode_timesteps,  

n_episodes, continue_training)  

Bases: tuple

continue_training  

    Alias for field number 3

episode_reward  

    Alias for field number 0

episode_timesteps  

    Alias for field number 1

n_episodes  

    Alias for field number 2

```

## 2.6.7 Values

```

class genrl.core.values.CnnCategoricalValue (*args, **kwargs)  

Bases: genrl.core.values.CnnNoisyValue  

Class for Categorical DQN's CNN Q-Value function

framestack  

    No. of frames being passed into the Q-value function

```

**Type** int  
**action\_dim**  
Action space dimensions

**Type** int  
**fc\_layers**  
Fully connected layer dimensions

**Type** tuple  
**noisy\_layers**  
Noisy layer dimensions

**Type** tuple  
**num\_atoms**  
Number of atoms used to discretise the Categorical DQN value distribution

**Type** int  
**forward**(*state*: *torch.Tensor*) → *torch.Tensor*  
Defines the computation performed at every call.

**Parameters** **state** (*Tensor*) – Input to value function

**class** *genrl.core.values.CnnDuelingValue*(\*args, \*\*kwargs)  
Bases: *genrl.core.values.CnnValue*

Class for Dueling DQN's MLP Q-Value function

**framestack**  
No. of frames being passed into the Q-value function

**Type** int  
**action\_dim**  
Action space dimensions

**Type** int  
**fc\_layers**  
Hidden layer dimensions

**Type** tuple  
**forward**(*inp*: *torch.Tensor*) → *torch.Tensor*  
Defines the computation performed at every call.

**Parameters** **state** (*Tensor*) – Input to value function

**class** *genrl.core.values.CnnNoisyValue*(\*args, \*\*kwargs)  
Bases: *genrl.core.values.CnnValue*, *genrl.core.values.MlpNoisyValue*

Class for Noisy DQN's CNN Q-Value function

**state\_dim**  
Number of previous frames to stack together

**Type** int  
**action\_dim**  
Action space dimensions

**Type** int

---

**fc\_layers**  
Fully connected layer dimensions  
**Type** tuple

**noisy\_layers**  
Noisy layer dimensions  
**Type** tuple

**num\_atoms**  
Number of atoms used to discretise the Categorical DQN value distribution  
**Type** int

**forward**(*state*: numpy.ndarray) → numpy.ndarray  
Defines the computation performed at every call.

**Parameters** **state** (*Tensor*) – Input to value function

**class** genrl.core.values.CnnValue(\*args, \*\*kwargs)  
Bases: *genrl.core.values.MlpValue*

CNN Value Function class

**param framestack** Number of previous frames to stack together  
**param action\_dim** Action dimension of environment  
**param val\_type** Specifies type of value function: ( “V” for V(s), “Qs” for Q(s), “Qsa” for Q(s,a))

**param fc\_layers** Sizes of hidden layers  
**type framestack** int  
**type action\_dim** int  
**type val\_type** string  
**type fc\_layers** tuple or list

**forward**(*state*: numpy.ndarray) → numpy.ndarray  
Defines the computation performed at every call.

**Parameters** **state** (*Tensor*) – Input to value function

**class** genrl.core.values.MlpCategoricalValue(\*args, \*\*kwargs)  
Bases: *genrl.core.values.MlpNoisyValue*

Class for Categorical DQN’s MLP Q-Value function

**state\_dim**  
Observation space dimensions  
**Type** int

**action\_dim**  
Action space dimensions  
**Type** int

**fc\_layers**  
Fully connected layer dimensions  
**Type** tuple

```
noisy_layers
    Noisy layer dimensions
        Type tuple

num_atoms
    Number of atoms used to discretise the Categorical DQN value distribution
        Type int

forward(state: torch.Tensor) → torch.Tensor
    Defines the computation performed at every call.

    Parameters state (Tensor) – Input to value function

class genrl.core.values.MlpDuelingValue (*args, **kwargs)
    Bases: genrl.core.values.MlpValue

    Class for Dueling DQN's MLP Q-Value function

    state_dim
        Observation space dimensions
            Type int

    action_dim
        Action space dimensions
            Type int

    hidden
        Hidden layer dimensions
            Type tuple

    forward(state: torch.Tensor) → torch.Tensor
        Defines the computation performed at every call.

        Parameters state (Tensor) – Input to value function

class genrl.core.values.MlpNoisyValue (*args, noisy_layers: Tuple = (128, 512), **kwargs)
    Bases: genrl.core.values.MlpValue

    reset_noise() → None
        Resets noise for any Noisy layers in Value function

class genrl.core.values.MlpValue (state_dim: int, action_dim: int = None, val_type: str = 'V',
                                    fc_layers: Tuple = (32, 32), **kwargs)
    Bases: genrl.core.base.BaseValue

    MLP Value Function class

        param state_dim State dimensions of environment
        param action_dim Action dimensions of environment
        param val_type Specifies type of value function: (
            "V" for V(s), "Qs" for Q(s), "Qsa" for Q(s,a))

            param hidden Sizes of hidden layers
            type state_dim int
            type action_dim int
            type val_type string
```

**type hidden** tuple or list

genrl.core.values.**get\_value\_from\_name** (*name*: str) → Union[Type[genrl.core.values.MlpValue], Type[genrl.core.values.CnnValue]]  
Gets the value function given the name of the value function

**Parameters** **name** (string) – Name of the value function needed

**Returns** Value function

## 2.7 Utilities

### 2.7.1 Logger

**class** genrl.utils.logger.**CSVLogger** (*logdir*: str)

Bases: object

CSV Logging class

**Parameters** **logdir** (string) – Directory to save log at

**close** () → None

Close the logger

**write** (*kvs*: Dict[str, Any], *log\_key*) → None

Add entry to logger

**Parameters** **kvs** (dict) – Entries to be logged

**class** genrl.utils.logger.**HumanOutputFormat** (*logdir*: str)

Bases: object

Output from a log file in a human readable format

**Parameters** **logdir** (string) – Directory at which log is present

**close** () → None

**max\_key\_len** (*kvs*: Dict[str, Any]) → None

Finds max key length

**Parameters** **kvs** (dict) – Entries to be logged

**round** (*num*: float) → float

Returns a rounded float value depending on self maxlen

**Parameters** **num** (float) – Value to round

**write** (*kvs*: Dict[str, Any], *log\_key*) → None

Log the entry out in human readable format

**Parameters** **kvs** (dict) – Entries to be logged

**write\_to\_file** (*kvs*: Dict[str, Any], *file*=<*\_io.TextIOWrapper* name='<stdout>' mode='w' encoding='UTF-8'>) → None

Log the entry out in human readable format

**Parameters**

- **kvs** (dict) – Entries to be logged

- **file** (*\_io.TextIOWrapper*) – Name of file to write logs to

```
class genrl.utils.logger.Logger(logdir: str = None, formats: List[str] = ['csv'])
```

Bases: object

Logger class to log important information

#### Parameters

- **logdir** (string) – Directory to save log at
- **formats** (list) – Formatting of each log ['csv', 'stdout', 'tensorboard']

**close**() → None

Close the logger

#### formats

Return save format(s)

#### logdir

Return log directory

**write**(kvs: Dict[str, Any], log\_key: str = 'timestep') → None

Add entry to logger

#### Parameters

- **kvs** (dict) – Entry to be logged
- **log\_key** (str) – Key plotted on log\_key

```
class genrl.utils.logger.TensorboardLogger(logdir: str)
```

Bases: object

Tensorboard Logging class

**Parameters** **logdir** (string) – Directory to save log at

**close**() → None

Close the logger

**write**(kvs: Dict[str, Any], log\_key: str = 'timestep') → None

Add entry to logger

#### Parameters

- **kvs** (dict) – Entries to be logged
- **log\_key** (str) – Key plotted on x\_axis

```
genrl.utils.logger.get_logger_by_name(name: str)
```

Gets the logger given the type of logger

**Parameters** **name** (string) – Name of the value function needed

**Returns** Logger

## 2.7.2 Utilities

```
genrl.utils.utils.cnn(channels: Tuple = (4, 16, 32), kernel_sizes: Tuple = (8, 4), strides: Tuple = (4, 2), **kwargs) → Tuple
```

(Generates a CNN model given input dimensions, channels, kernel\_sizes and  
strides)

**param channels** Input output channels before and after each convolution

---

**param kernel\_sizes** Kernel sizes for each convolution

**param strides** Strides for each convolution

**param in\_size** Input dimensions (assuming square input)

**type channels** tuple

**type kernel\_sizes** tuple

**type strides** tuple

**type in\_size** int

**returns** (Convolutional Neural Network with convolutional layers and activation layers)

```
genrl.utils.utils.get_env_properties (env: Union[gym.core.Env,
                                                genrl.environments.vec_env.vector_envs.VecEnv],
                                                network: Union[str, Any] = 'mlp') → Tuple[int]
```

Finds important properties of environment

**param env** Environment that the agent is interacting with

**type env** Gym Environment

**param network** Type of network architecture, eg. “mlp”, “cnn”

**type network** str

**returns** (State space dimensions, Action space dimensions,

**discreteness of action space and action limit (highest action value)**

**rtype** int, float, ... ; int, float, ... ; bool; int, float, ...

```
genrl.utils.utils.get_model (type_: str, name_: str) → Union
```

Utility to get the class of required function

**param type\_** “ac” for Actor Critic, “v” for Value, “p” for Policy

**param name\_** Name of the specific structure of model. (

Eg. “mlp” or “cnn”)

**type type\_** string

**returns** Required class. Eg. MlpActorCritic

```
genrl.utils.utils.mlp (sizes: Tuple, activation: str = 'relu', sac: bool = False)
```

Generates an MLP model given sizes of each layer

**param sizes** Sizes of hidden layers

**param sac** True if Soft Actor Critic is being used, else False

**type sizes** tuple or list

**type sac** bool

**returns** (Neural Network with fully-connected linear layers and activation layers)

```
genrl.utils.utils.noisy_mlp(fc_layers: List[int], noisy_layers: List[int], activation='relu')  
    Noisy MLP generating helper function
```

**Parameters**

- **fc\_layers** (list of int) – List of fully connected layers
- **noisy\_layers** (list of int) – List of noisy layers
- **activation** (str) – Activation function to be used. [“tanh”, “relu”]

**Returns** Noisy MLP model

```
genrl.utils.utils.safe_mean(log: List[int])
```

Returns 0 if there are no elements in logs

```
genrl.utils.utils.set_seeds(seed: int, env: Union[gym.core.Env,  
                                                genrl.environments.vec_env.vector_envs.VecEnv] = None) →  
    None
```

Sets seeds for reproducibility

**Parameters**

- **seed** (int) – Seed Value
- **env** (Gym Environment) – Optionally pass gym environment to set its seed

## 2.7.3 Models

```
class genrl.utils.models.TabularModel(s_dim: int, a_dim: int)  
Bases: object
```

Sample-based tabular model class for deterministic, discrete environments

**Parameters**

- **s\_dim** (int) – environment state dimension
- **a\_dim** (int) – environment action dimension

```
add(state: numpy.ndarray, action: numpy.ndarray, reward: float, next_state: numpy.ndarray) → None  
    add transition to model :param state: state :param action: action :param reward: reward :param next_state:  
    next state :type state: float array :type action: int :type reward: int :type next_state: float array
```

```
is_empty() → bool
```

Check if the model has been updated or not

**Returns** True if model not updated yet**Return type** bool

```
sample() → Tuple
```

sample state action pair from model

**Returns** state and action**Return type** int, float, .. ; int, float, ..

```
step(state: numpy.ndarray, action: numpy.ndarray) → Tuple
```

return consequence of action at state

**Returns** reward and next state**Return type** int; int, float, ..

---

```
genrl.utils.models.get_model_from_name(name_: str)
    get model object from name
```

**Parameters** `name` (`str`) – name of the model [‘tabular’]

**Returns** the model

## 2.8 Trainers

### 2.8.1 On-Policy Trainer

On Policy Trainer Class

Trainer class for all the On Policy Agents: A2C, PPO1 and VPG

```
genrl.trainers.OnPolicyTrainer.agent
```

Agent algorithm object

**Type** object

```
genrl.trainers.OnPolicyTrainer.env
```

Environment

**Type** object

```
genrl.trainers.OnPolicyTrainer.log_mode
```

List of different kinds of logging. Supported: [“csv”, “stdout”, “tensorboard”]

**Type** list of str

```
genrl.trainers.OnPolicyTrainer.log_key
```

Key plotted on x\_axis. Supported: [“timestep”, “episode”]

**Type** str

```
genrl.trainers.OnPolicyTrainer.log_interval
```

Timesteps between successive logging of parameters onto the console

**Type** int

```
genrl.trainers.OnPolicyTrainer.logdir
```

Directory where log files should be saved.

**Type** str

```
genrl.trainers.OnPolicyTrainer.epochs
```

Total number of epochs to train for

**Type** int

```
genrl.trainers.OnPolicyTrainer.max_timesteps
```

Maximum limit of timesteps to train for

**Type** int

```
genrl.trainers.OnPolicyTrainer.off_policy
```

True if the agent is an off policy agent, False if it is on policy

**Type** bool

```
genrl.trainers.OnPolicyTrainer.save_interval
```

Timesteps between successive saves of the agent’s important hyperparameters

**Type** int

genrl.trainers.OnPolicyTrainer.**save\_model**

Directory where the checkpoints of agent parameters should be saved

**Type** str

genrl.trainers.OnPolicyTrainer.**run\_num**

A run number allotted to the save of parameters

**Type** int

genrl.trainers.OnPolicyTrainer.**load\_model**

File to load saved parameter checkpoint from

**Type** str

genrl.trainers.OnPolicyTrainer.**render**

True if environment is to be rendered during training, else False

**Type** bool

genrl.trainers.OnPolicyTrainer.**evaluate\_episodes**

Number of episodes to evaluate for

**Type** int

genrl.trainers.OnPolicyTrainer.**seed**

Set seed for reproducibility

**Type** int

genrl.trainers.OnPolicyTrainer.**n\_envs**

Number of environments

## 2.8.2 Off-Policy Trainer

Off Policy Trainer Class

Trainer class for all the Off Policy Agents: DQN (all variants), DDPG, TD3 and SAC

genrl.trainers.OffPolicyTrainer.**agent**

Agent algorithm object

**Type** object

genrl.trainers.OffPolicyTrainer.**env**

Environment

**Type** object

genrl.trainers.OffPolicyTrainer.**buffer**

Replay Buffer object

**Type** object

genrl.trainers.OffPolicyTrainer.**max\_ep\_len**

Maximum Episode length for training

**Type** int

genrl.trainers.OffPolicyTrainer.**warmup\_steps**

Number of warmup steps. (random actions are taken to add randomness to training)

**Type** int

---

genrl.trainers.OffPolicyTrainer.**start\_update**  
 Timesteps after which the agent networks should start updating  
**Type** int

genrl.trainers.OffPolicyTrainer.**update\_interval**  
 Timesteps between target network updates  
**Type** int

genrl.trainers.OffPolicyTrainer.**log\_mode**  
 List of different kinds of logging. Supported: [“csv”, “stdout”, “tensorboard”]  
**Type** list of str

genrl.trainers.OffPolicyTrainer.**log\_key**  
 Key plotted on x\_axis. Supported: [“timestep”, “episode”]  
**Type** str

genrl.trainers.OffPolicyTrainer.**log\_interval**  
 Timesteps between successive logging of parameters onto the console  
**Type** int

genrl.trainers.OffPolicyTrainer.**logdir**  
 Directory where log files should be saved.  
**Type** str

genrl.trainers.OffPolicyTrainer.**epochs**  
 Total number of epochs to train for  
**Type** int

genrl.trainers.OffPolicyTrainer.**max\_timesteps**  
 Maximum limit of timesteps to train for  
**Type** int

genrl.trainers.OffPolicyTrainer.**off\_policy**  
 True if the agent is an off policy agent, False if it is on policy  
**Type** bool

genrl.trainers.OffPolicyTrainer.**save\_interval**  
 Timesteps between successive saves of the agent’s important hyperparameters  
**Type** int

genrl.trainers.OffPolicyTrainer.**save\_model**  
 Directory where the checkpoints of agent parameters should be saved  
**Type** str

genrl.trainers.OffPolicyTrainer.**rung\_num**  
 A run number allotted to the save of parameters  
**Type** int

genrl.trainers.OffPolicyTrainer.**load\_model**  
 File to load saved parameter checkpoint from  
**Type** str

genrl.trainers.OffPolicyTrainer.**render**  
 True if environment is to be rendered during training, else False

**Type** bool  
genrl.trainers.OffPolicyTrainer.**evaluate\_episodes**  
Number of episodes to evaluate for  
**Type** int  
genrl.trainers.OffPolicyTrainer.**seed**  
Set seed for reproducibility  
**Type** int  
genrl.trainers.OffPolicyTrainer.**n\_envs**  
Number of environments

### 2.8.3 Classical Trainer

Global trainer class for classical RL algorithms

**param agent** Algorithm object to train  
**param env** standard gym environment to train on  
**param mode** mode of value function update ['learn', 'plan', 'dyna']  
**param model** model to use for planning ['tabular']  
**param n\_episodes** number of training episodes  
**param plan\_n\_steps** number of planning step per environment interaction  
**param start\_steps** number of initial exploration timesteps  
**param seed** seed for random number generator  
**param render** render gym environment  
**type agent** object  
**type env** Gym environment  
**type mode** str  
**type model** str  
**type n\_episodes** int  
**type plan\_n\_steps** int  
**type start\_steps** int  
**type seed** int  
**type render** bool

### 2.8.4 Deep Contextual Bandit Trainer

Bandit Trainer Class

**param agent** Agent to train.  
**type agent** genrl.deep.bandit.dcb\_agents.DCBAgent  
**param bandit** Bandit to train agent on.  
**type bandit** genrl.deep.bandit.data\_bandits.DataBasedBandit

---

**param logdir** Path to directory to store logs in.  
**type logdir** str  
**param log\_mode** List of modes for logging.  
**type log\_mode** List[str]

## 2.8.5 Multi Armed Bandit Trainer

Bandit Trainer Class

**param agent** Agent to train.  
**type agent** genrl.deep.bandit.dcb\_agents.DCBAgent  
**param bandit** Bandit to train agent on.  
**type bandit** genrl.deep.bandit.data\_bandits.DataBasedBandit  
**param logdir** Path to directory to store logs in.  
**type logdir** str  
**param log\_mode** List of modes for logging.  
**type log\_mode** List[str]

## 2.8.6 Base Trainer

Base Trainer Class

To be inherited specific use-cases

genrl.trainers.Trainer.**agent**  
Agent algorithm object  
**Type** object  
genrl.trainers.Trainer.**env**  
Environment  
**Type** object  
genrl.trainers.Trainer.**log\_mode**  
List of different kinds of logging. Supported: [“csv”, “stdout”, “tensorboard”]  
**Type** list of str  
genrl.trainers.Trainer.**log\_key**  
Key plotted on x\_axis. Supported: [“timestep”, “episode”]  
**Type** str  
genrl.trainers.Trainer.**log\_interval**  
Timesteps between successive logging of parameters onto the console  
**Type** int  
genrl.trainers.Trainer.**logdir**  
Directory where log files should be saved.  
**Type** str

```
genrl.trainers.Trainer.epochs
    Total number of epochs to train for
        Type int

genrl.trainers.Trainer.max_timesteps
    Maximum limit of timesteps to train for
        Type int

genrl.trainers.Trainer.off_policy
    True if the agent is an off policy agent, False if it is on policy
        Type bool

genrl.trainers.Trainer.save_interval
    Timesteps between successive saves of the agent's important hyperparameters
        Type int

genrl.trainers.Trainer.save_model
    Directory where the checkpoints of agent parameters should be saved
        Type str

genrl.trainers.Trainer.run_num
    A run number allotted to the save of parameters
        Type int

genrl.trainers.Trainer.load_model
    File to load saved parameter checkpoint from
        Type str

genrl.trainers.Trainer.render
    True if environment is to be rendered during training, else False
        Type bool

genrl.trainers.Trainer.evaluate_episodes
    Number of episodes to evaluate for
        Type int

genrl.trainers.Trainer.seed
    Set seed for reproducibility
        Type int

genrl.trainers.Trainer.n_envs
    Number of environments
```

## 2.9 Common

### 2.9.1 Classical Common

```
genrl.classical.common.models
genrl.classical.common.trainer
```

**genrl.classical.common.values****2.9.2 Bandit Common****genrl.bandit.core****genrl.bandit.trainer****genrl.bandit.agents.cb\_agents.common.base\_model**

**class** genrl.agents.bandits.contextual.common.base\_model.**Model** (*layer*, *\*\*kwargs*)

Bases: torch.nn.modules.module.Module, abc.ABC

Bayesian Neural Network used in Deep Contextual Bandit Models.

**Parameters**

- **context\_dim** (*int*) – Length of context vector.
- **hidden\_dims** (*List[int]*, *optional*) – Dimensions of hidden layers of network.
- **n\_actions** (*int*) – Number of actions that can be selected. Taken as length of output vector for network to predict.
- **init\_lr** (*float*, *optional*) – Initial learning rate.
- **max\_grad\_norm** (*float*, *optional*) – Maximum norm of gradients for gradient clipping.
- **lr\_decay** (*float*, *optional*) – Decay rate for learning rate.
- **lr\_reset** (*bool*, *optional*) – Whether to reset learning rate ever train interval. Defaults to False.
- **dropout\_p** (*Optional[float]*, *optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **noise\_std** (*float*) – Standard deviation of noise used in the network. Defaults to 0.1

**use\_dropout**

Indicated whether or not dropout should be used in forward pass.

**Type** int

**forward** (*context*: torch.Tensor, *\*\*kwargs*) → Dict[str, torch.Tensor]

Computes forward pass through the network.

**Parameters** **context** (*torch.Tensor*) – The context vector to perform forward pass on.

**Returns** Dictionary of outputs

**Return type** Dict[str, torch.Tensor]

**train\_model** (*db*: genrl.agents.bandits.contextual.common.transition.TransitionDB, *epochs*: *int*, *batch\_size*: *int*)

Trains the network on a given database for given epochs and batch\_size.

**Parameters**

- **db** (*TransitionDB*) – The database of transitions to train on.
- **epochs** (*int*) – Number of gradient steps to take.
- **batch\_size** (*int*) – The size of each batch to perform gradient descent on.

**genrl.bandit.agents.cb\_agents.common.bayesian**

```
class genrl.agents.bandits.contextual.common.bayesian.BayesianLinear(in_features:  
    int,  
    out_features:  
    int,  
    bias:  
    bool =  
    True)
```

Bases: torch.nn.modules.module.Module

Linear Layer for Bayesian Neural Networks.

**Parameters**

- **in\_features** (*int*) – size of each input sample
- **out\_features** (*int*) – size of each output sample
- **bias** (*bool, optional*) – Whether to use an additive bias. Defaults to True.

**forward**(*x: torch.Tensor, kl: bool = True, frozen: bool = False*) → Tuple[*torch.Tensor, Optional[torch.Tensor]*]  
Apply linear transformation to input.

The weight and bias is sampled for each forward pass from a normal distribution. The KL divergence of the sampled weight and bias can also be computed if specified.

**Parameters**

- **x** (*torch.Tensor*) – Input to be transformed
- **kl** (*bool, optional*) – Whether to compute the KL divergence. Defaults to True.
- **frozen** (*bool, optional*) – Whether to freeze current parameters. Defaults to False.

**Returns**

The transformed input and optionally the computed KL divergence value.

**Return type** Tuple[*torch.Tensor, Optional[torch.Tensor]*]

**reset\_parameters**() → None

Resets weight and bias parameters of the layer.

```
class genrl.agents.bandits.contextual.common.bayesian.BayesianNNBanditModel(**kwargs)
```

Bases: *genrl.agents.bandits.contextual.common.base\_model.Model*

Bayesian Neural Network used in Deep Contextual Bandit Models.

**Parameters**

- **context\_dim** (*int*) – Length of context vector.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network.
- **n\_actions** (*int*) – Number of actions that can be selected. Taken as length of output vector for network to predict.
- **init\_lr** (*float, optional*) – Initial learning rate.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate.

- **lr\_reset** (*bool, optional*) – Whether to reset learning rate ever train interval. Defaults to False.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.
- **noise\_std** (*float*) – Standard deviation of noise used in the network. Defaults to 0.1

**use\_dropout**

Indicated whether or not dropout should be used in forward pass.

**Type** int

**forward** (*context: torch.Tensor, kl: bool = True*) → Dict[str, torch.Tensor]

Computes forward pass through the network.

**Parameters** **context** (*torch.Tensor*) – The context vector to perform forward pass on.

**Returns** Dictionary of outputs

**Return type** Dict[str, torch.Tensor]

**genrl.bandit.agents.cb\_agents.common.neural**

```
class genrl.agents.bandits.contextual.common.neural.NeuralBanditModel(**kwargs)
Bases: genrl.agents.bandits.contextual.common.base_model.Model
```

Neural Network used in Deep Contextual Bandit Models.

**Parameters**

- **context\_dim** (*int*) – Length of context vector.
- **hidden\_dims** (*List[int], optional*) – Dimensions of hidden layers of network.
- **n\_actions** (*int*) – Number of actions that can be selected. Taken as length of output vector for network to predict.
- **init\_lr** (*float, optional*) – Initial learning rate.
- **max\_grad\_norm** (*float, optional*) – Maximum norm of gradients for gradient clipping.
- **lr\_decay** (*float, optional*) – Decay rate for learning rate.
- **lr\_reset** (*bool, optional*) – Whether to reset learning rate ever train interval. Defaults to False.
- **dropout\_p** (*Optional[float], optional*) – Probability for dropout. Defaults to None which implies dropout is not to be used.

**use\_dropout**

Indicated whether or not dropout should be used in forward pass.

**Type** bool

**forward** (*context: torch.Tensor*) → Dict[str, torch.Tensor]

Computes forward pass through the network.

**Parameters** **context** (*torch.Tensor*) – The context vector to perform forward pass on.

**Returns** Dictionary of outputs

**Return type** Dict[str, torch.Tensor]

**genrl.bandit.agents.cb\_agents.common.transition**

```
class genrl.agents.bandits.contextual.common.transition.TransitionDB(device:  
    Union[str,  
          torch.device]  
    =  
    'cpu')
```

Bases: object

Database for storing (context, action, reward) transitions.

**Parameters** **device** (*str*) – Device to use for tensor operations. “cpu” for cpu or “cuda” for cuda.  
Defaults to “cpu”.

**db**

Dictionary containing list of transitions.

**Type** dict

**db\_size**

Number of transitions stored in database.

**Type** int

**device**

Device to use for tensor operations.

**Type** torch.device

**add** (*context: torch.Tensor, action: int, reward: int*)  
Add (context, action, reward) transition to database

**Parameters**

- **context** (*torch.Tensor*) – Context received
- **action** (*int*) – Action taken
- **reward** (*int*) – Reward received

**get\_data** (*batch\_size: Optional[int] = None*) → Tuple[*torch.Tensor, torch.Tensor, torch.Tensor*]  
Get a batch of transition from database

**Parameters** **batch\_size** (*Union[int, None], optional*) – Size of batch required.  
Defaults to None which implies all transitions in the database are to be included in batch.

**Returns**

**Tuple of stacked** contexts, actions, rewards tensors.

**Return type** Tuple[*torch.Tensor, torch.Tensor, torch.Tensor*]

**get\_data\_for\_action** (*action: int, batch\_size: Optional[int] = None*) → Tuple[*torch.Tensor, torch.Tensor*]  
Get a batch of transition from database for a given action.

**Parameters**

- **action** (*int*) – The action to sample transitions for.
- **batch\_size** (*Union[int, None], optional*) – Size of batch required. Defaults to None which implies all transitions in the database are to be included in batch.

**Returns**

**Tuple of stacked** contexts and rewards tensors.

**Return type** Tuple[torch.Tensor, torch.Tensor]



---

## Python Module Index

---

**g**

genrl.agents.bandits.multiarmed.ucb, 72  
genrl.agents.bandits.contextual.base, 57  
genrl.agents.bandits.contextual.bootstrap\_neural, 58  
genrl.agents.bandits.contextual.common.base\_model, 107  
genrl.agents.bandits.contextual.common.bayesian, 108  
genrl.agents.bandits.contextual.common.noisy, 109  
genrl.agents.bandits.contextual.common.transition, 110  
genrl.agents.bandits.contextual.fixed, 59  
genrl.agents.bandits.contextual.linpos, 60  
genrl.agents.bandits.contextual.neural\_greedy, 61  
genrl.agents.bandits.contextual.neural\_laplace, 62  
genrl.agents.bandits.contextual.neural\_nsgail, 63  
genrl.agents.bandits.contextual.variation, 65  
genrl.agents.bandits.multiarmed.base, 66  
genrl.agents.bandits.multiarmed.bayesian, 67  
genrl.agents.bandits.multiarmed.bernoulli, 68  
genrl.agents.bandits.multiarmed.epsgreedy, 69  
genrl.agents.bandits.multiarmed.gaussian, 70  
genrl.agents.bandits.multiarmed.gradient, 70  
genrl.agents.bandits.multiarmed.thompson, 71  
genrl.agents.classical.qlearning.qlearning, 55  
genrl.agents.classical.sarsa.sarsa, 56  
genrl.agents.deep.a2c.a2c, 31  
genrl.agents.deep.ddpg.ddpg, 33  
genrl.agents.deep.dqn.base, 35  
genrl.agents.deep.dqn.categorical, 38  
genrl.agents.deep.dqn.double, 40  
genrl.agents.deep.dqn.dueling, 41  
genrl.agents.deep.dqn.noisy, 43  
genrl.agents.deep.dqn.prioritized, 44  
genrl.agents.deep.dqn.utils, 45  
genrl.agents.deep.ppo1.ppo1, 47  
genrl.agents.deep.sac.sac, 53  
genrl.agents.deep.td3.td3, 51  
genrl.agents.deep.vpg.vpg, 49  
genrl.core.actor\_critic, 82  
genrl.core.base, 85  
genrl.core.buffers, 87  
genrl.core.noise, 89  
genrl.core.rollout\_storage, 91  
genrl.core.values, 93  
genrl.environments, 82  
genrl.environments.action\_wrappers, 76  
genrl.environments.atari\_preprocessing, 76  
genrl.environments.atari\_wrappers, 78  
genrl.environments.base\_wrapper, 78  
genrl.environments.frame\_stack, 79  
genrl.environments.gym\_wrapper, 80  
genrl.environments(suite, 80  
genrl.environments.time\_limit, 81  
genrl.environments.vec\_env, 76  
genrl.environments.vec\_env.monitor, 73  
genrl.environments.vec\_env.normalize, 73  
genrl.environments.vec\_env.utils, 74

```
genrl.environments.vec_env.vector_envs,  
    74  
genrl.environments.vec_env.wrappers, 76  
genrl.trainers.ClassicalTrainer, 104  
genrl.trainers.DCBTrainer, 104  
genrl.trainers.MABTrainer, 105  
genrl.trainers.OffPolicyTrainer, 102  
genrl.trainers.OnPolicyTrainer, 101  
genrl.trainers.Trainer, 105  
genrl.utils.logger, 97  
genrl.utils.models, 100  
genrl.utils.utils, 98
```

---

## Index

---

### A

a (*genrl.agents.bandits.multiarmed.bayesian.BayesianUCBMABAgen* attribute), 93  
activation (*genrl.core.actor\_critic.MlpActorCritic* attribute), 67  
a (*genrl.agents.bandits.multiarmed.thompson.ThompsonSamplingMABAgen* attribute), 84  
activation (*genrl.core.actor\_critic.MlpSingleActorMultiCritic* attribute), 71  
attribute), 84  
A2C (*class in genrl.agents.deep.a2c.a2c*), 31  
action () (*genrl.environments.action\_wrappers.ClipAction* method), 76  
action () (*genrl.environments.action\_wrappers.RescaleAction* method), 77  
action\_dim (*genrl.core.actor\_critic.MlpActorCritic* attribute), 83  
action\_dim (*genrl.core.actor\_critic.MlpSingleActorMultiCritic* attribute), 84  
action\_dim (*genrl.core.values.CnnCategoricalValue* attribute), 94  
action\_dim (*genrl.core.values.CnnDuelingValue* attribute), 94  
action\_dim (*genrl.core.values.CnnNoisyValue* attribute), 94  
action\_dim (*genrl.core.values.MlpCategoricalValue* attribute), 95  
action\_dim (*genrl.core.values.MlpDuelingValue* attribute), 96  
action\_hist (*genrl.agents.bandits.multiarmed.base.MABAgen* attribute), 66  
action\_shape (*genrl.environments.gym\_wrapper.GymWrapper* attribute), 80  
action\_shape (*genrl.environments.vec\_env.vector\_envs.VecEnv* attribute), 75  
action\_spaces (*genrl.environments.vec\_env.vector\_envs.VecEnv* attribute), 75  
ActionNoise (*class in genrl.core.noise*), 89  
actions (*genrl.core.buffers.PrioritizedReplayBufferSamples* attribute), 88  
actions (*genrl.core.buffers.ReplayBufferSamples* attribute), 88  
actions (*genrl.core.rollout\_storage.ReplayBufferSamples* attribute), 92  
actions (*genrl.core.rollout\_storage.RolloutBufferSamples* attribute), 93  
activation (*genrl.core.actor\_critic.MlpActorCritic* attribute), 84  
add () (*genrl.agents.bandits.contextual.common.transition.TransitionDB* method), 110  
add () (*genrl.core.rollout\_storage.BaseBuffer* method), 91  
add () (*genrl.core.rollout\_storage.RolloutBuffer* method), 92  
add () (*genrl.utils.models.TabularModel* method), 100  
advantages (*genrl.core.rollout\_storage.RolloutBufferSamples* attribute), 93  
agent (*in module genrl.trainers.OffPolicyTrainer*), 102  
agent (*in module genrl.trainers.OnPolicyTrainer*), 101  
agent (*in module genrl.trainers.Trainer*), 105  
alpha (*genrl.agents.bandits.multiarmed.gradient.GradientMABAgen* attribute), 70  
alpha (*genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN* attribute), 45  
alpha (*genrl.agents.deep.sac.sac.SAC* attribute), 54  
AtariEnv () (*in module genrl.environments.suite*), 80  
AtariPreprocessing (*class in genrl.environments.atari\_preprocessing*), 77  
AtariTimeLimit (*class in genrl.environments.time\_limit*), 81  
B

### B

b (*genrl.agents.bandits.multiarmed.bayesian.BayesianUCBMABAgen* attribute), 67  
bandit (*genrl.agents.bandits.contextual.base.DCBAgent* attribute), 57  
BaseActorCritic (*class in genrl.core.base*), 85  
BaseBuffer (*class in genrl.core.rollout\_storage*), 91  
BasePolicy (*class in genrl.core.base*), 86

BaseValue (class in `genrl.core.base`), 86  
 BaseWrapper (class in `genrl.environments.base_wrapper`), 78  
 batch\_size (`genrl.agents.deep.a2c.a2c.A2C` attribute), 31  
 batch\_size (`genrl.agents.deep.ddpg.ddpg.DDPG` attribute), 33  
 batch\_size (`genrl.agents.deep.dqn.base.DQN` attribute), 35  
 batch\_size (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
 batch\_size (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 40  
 batch\_size (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 42  
 batch\_size (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
 batch\_size (`genrl.agents.deep.ppo1.ppo1.PPO1` attribute), 48  
 batch\_size (`genrl.agents.deep.sac.sac.SAC` attribute), 54  
 batch\_size (`genrl.agents.deep.td3.td3.TD3` attribute), 52

**C**

calculate\_epsilon\_by\_frame ()  
 categorical\_greedy\_action () (in module `genrl.agents.deep.dqn.utils`), 45  
 categorical\_q\_loss () (in module `genrl.agents.deep.dqn.utils`), 46  
 categorical\_q\_target () (in module `genrl.agents.deep.dqn.utils`), 46  
 categorical\_q\_values () (in module `genrl.agents.deep.dqn.utils`), 46

clip\_param (`genrl.agents.deep.ppo1.ppo1.PPO1` attribute), 48

close () (class in `genrl.environments.action_wrappers`), 76  
 close () (genrl.environments.base\_wrapper.BaseWrapper method), 78  
 close () (genrl.environments.gym\_wrapper.GymWrapper method), 80  
 close () (genrl.environments.vec\_env.normalize.VecNormalize method), 74

close () (genrl.environments.vec\_env.vector\_envs.SerialVecEnv method), 75  
 close () (genrl.environments.vec\_env.vector\_envs.SubProcessVecEnv method), 75  
 close () (genrl.environments.vec\_env.vector\_envs.VecEnv method), 75  
 close () (genrl.environments.vec\_env.wrappers.VecEnvWrapper method), 76  
 close () (genrl.utils.logger.CSVLogger method), 97  
 close () (genrl.utils.logger.HumanOutputFormat method), 97  
 close () (genrl.utils.logger.Logger method), 98  
 close () (genrl.utils.logger.TensorboardLogger method), 98  
 cnn () (in module `genrl.utils.utils`), 98

CNNActorCritic (*class in genrl.core.actor\_critic*), 82  
 CnnCategoricalValue (*class in genrl.core.values*), 93  
 CnnDuelingValue (*class in genrl.core.values*), 94  
 CnnNoisyValue (*class in genrl.core.values*), 94  
 CNNPolicy (*class in genrl.core.policies*), 90  
 CnnValue (*class in genrl.core.values*), 95  
 compute\_returns\_and\_advantage()  
     (*genrl.core.rollout\_storage.RolloutBuffer method*), 93  
 confidence (*genrl.agents.bandits.multiarmed.bayesian.BayesianUCBMAgent* and *deep.ppo1.PPO1 attribute*), 48  
 confidence (*genrl.agents.bandits.multiarmed.ucb.UCBMAgent attribute*), 72  
 continue\_training  
     (*genrl.core.rollout\_storage.RolloutReturn attribute*), 93  
 counts (*genrl.agents.bandits.multiarmed.base.MABAgent attribute*), 66  
 create\_model  
     (*genrl.agents.deep.a2c.a2c.A2C attribute*), 31  
 create\_model  
     (*genrl.agents.deep.ddpg.ddpg.DDPG attribute*), 33  
 create\_model  
     (*genrl.agents.deep.dqn.base.DQN attribute*), 35  
 create\_model  
     (*genrl.agents.deep.dqn.categorical.CategoricalDQN attribute*), 38  
 create\_model  
     (*genrl.agents.deep.ppo1.ppo1.PPO1 attribute*), 47  
 create\_model  
     (*genrl.agents.deep.sac.sac.SAC attribute*), 53  
 create\_model  
     (*genrl.agents.deep.td3.td3.TD3 attribute*), 51  
 CSVLogger (*class in genrl.utils.logger*), 97

**D**

db (*genrl.agents.bandits.contextual.common.transition.TransitionDB attribute*), 110  
 db\_size (*genrl.agents.bandits.contextual.common.transition.TransitionDB attribute*), 110  
 DCBAGent (*class in genrl.agents.bandits.contextual.base*), 57  
 DDPG (*class in genrl.agents.deep.ddpg.ddpg*), 33  
 ddqn\_q\_target()  
     (*in module genrl.agents.deep.dqn.utils*), 46  
 device (*genrl.agents.bandits.contextual.base.DCBAGent attribute*), 57  
 device (*genrl.agents.bandits.contextual.common.transition.TransitionDB attribute*), 110  
 device (*genrl.agents.deep.a2c.a2c.A2C attribute*), 32  
 device (*genrl.agents.deep.ddpg.ddpg.DDPG attribute*), 34  
 device (*genrl.agents.deep.dqn.base.DQN attribute*), 36  
 device (*genrl.agents.deep.dqn.categorical.CategoricalDQN attribute*), 39  
 device (*genrl.agents.deep.dqn.double.DoubleDQN attribute*), 41  
 device (*genrl.agents.deep.dqn.dueling.DuelingDQN attribute*), 42  
 device (*genrl.agents.deep.dqn.noisy.NoisyDQN attribute*), 44  
 device (*genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN attribute*), 45  
 device (*genrl.agents.deep.sac.sac.SAC attribute*), 54  
 device (*genrl.agents.deep.td3.td3.TD3 attribute*), 52  
 discrete (*genrl.core.actor\_critic.MlpActorCritic attribute*), 83  
 discrete (*genrl.core.actor\_critic.MlpSingleActorMultiCritic attribute*), 84  
 dones (*genrl.core.buffers.PrioritizedReplayBufferSamples attribute*), 88  
 dones (*genrl.core.buffers.ReplayBufferSamples attribute*), 88  
 dones (*genrl.core.rollout\_storage.ReplayBufferSamples attribute*), 92  
 DoubleDQN (*class in genrl.agents.deep.dqn.double*), 40  
 DoubleDQN (*in genrl.agents.deep.dqn.base*), 35  
 DuelingDQN (*class in genrl.agents.deep.dqn.dueling*), 41

**E**

empty\_logs()  
     (*genrl.agents.deep.a2c.a2c.A2C method*), 32  
 empty\_logs()  
     (*genrl.agents.deep.ddpg.ddpg.DDPG method*), 34  
 empty\_logs()  
     (*genrl.agents.deep.dqn.base.DQN method*), 36  
 empty\_logs()  
     (*genrl.agents.deep.ppo1.ppo1.PPO1 method*), 48  
 empty\_logs()  
     (*genrl.agents.deep.sac.sac.SAC method*), 54  
 empty\_logs()  
     (*genrl.agents.deep.td3.td3.TD3 method*), 52  
 empty\_logs()  
     (*genrl.agents.deep.vpg.vpg.VPG method*), 50  
 entropy\_coeff  
     (*genrl.agents.deep.a2c.a2c.A2C attribute*), 32  
 entropy\_coeff  
     (*genrl.agents.deep.ppo1.ppo1.PPO1 attribute*), 48  
 entropy\_tuning  
     (*genrl.agents.deep.sac.sac.SAC attribute*), 54  
 env (*genrl.agents.classical.qlearning.QLearning attribute*), 56  
 env  
     (*genrl.agents.classical.sarsa.sarsa.SARSA attribute*), 56

env (`genrl.agents.deep.a2c.a2c.A2C` attribute), 31  
 env (`genrl.agents.deep.ddpg.ddpg.DDPG` attribute), 33  
 env (`genrl.agents.deep.dqn.base.DQN` attribute), 35  
 env (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
 env (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 40  
 env (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 41  
 env (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
 env (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN` attribute), 44  
 env (`genrl.agents.deep.ppo1.ppo1.PPO1` attribute), 47  
 env (`genrl.agents.deep.sac.sac.SAC` attribute), 53  
 env (`genrl.agents.deep.td3.td3.TD3` attribute), 51  
 env (in module `genrl.trainers.OffPolicyTrainer`), 102  
 env (in module `genrl.trainers.OnPolicyTrainer`), 101  
 env (in module `genrl.trainers.Trainer`), 105  
 episode\_reward (`genrl.core.rollout_storage.RolloutReturn` attribute), 93  
 episode\_timesteps (`genrl.core.rollout_storage.RolloutReturn` attribute), 93  
 epochs (in module `genrl.trainers.OffPolicyTrainer`), 103  
 epochs (in module `genrl.trainers.OnPolicyTrainer`), 101  
 epochs (in module `genrl.trainers.Trainer`), 105  
 eps (`genrl.agents.bandits.multiarmed.EpsGreedy` attribute), 69  
 EpsGreedyMABAgent (class in `genrl.agents.bandits.multiarmed.epsgreedy`), 69  
 epsilon (`genrl.agents.classical.qlearning.QLearning` attribute), 56  
 epsilon (`genrl.agents.classical.sarsa.sarsa.SARSA` attribute), 56  
 epsilon\_decay (`genrl.agents.deep.dqn.base.DQN` attribute), 36  
 epsilon\_decay (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
 epsilon\_decay (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 41  
 epsilon\_decay (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 42  
 epsilon\_decay (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
 epsilon\_decay (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN` attribute), 45  
 evaluate\_actions () (`genrl.agents.deep.a2c.a2c.A2C` method), 32  
 evaluate\_actions () (`genrl.agents.deep.ppo1.ppo1.PPO1` method), 48  
 evaluate\_episodes (in module `genrl.trainers.OffPolicyTrainer`), 104  
 evaluate\_episodes (in module `genrl.trainers.OnPolicyTrainer`), 102  
 evaluate\_episodes (in module `genrl.trainers.Trainer`), 106  
 extend () (`genrl.core.buffers.ReplayBuffer` method), 88  
 extend () (`genrl.core.rollout_storage.BaseBuffer` method), 91

**F**

fc\_layers (`genrl.core.values.CnnCategoricalValue` attribute), 94  
 fc\_layers (`genrl.core.values.CnnDuelingValue` attribute), 94  
 fc\_layers (`genrl.core.values.CnnNoisyValue` attribute), 94  
 fc\_layers (`genrl.core.values.MlpCategoricalValue` attribute), 95  
 FireReset (class in `genrl.environments.atari_wrappers`), 78  
 FixedAgent (class in `genrl.agents.bandits.contextual.fixed`), 59  
 formats (`genrl.utils.logger.Logger` attribute), 98  
 forward () (`genrl.agents.bandits.contextual.common.base_model.Model` method), 107  
 forward () (`genrl.agents.bandits.contextual.common.bayesian.BayesianL` method), 109  
 forward () (`genrl.agents.bandits.contextual.common.bayesian.BayesianN` method), 109  
 forward () (`genrl.agents.bandits.contextual.common.neural.NeuralBandit` method), 109  
 forward () (`genrl.core.actor_critic.MlpSingleActorMultiCritic` method), 84  
 forward () (`genrl.core.base.BasePolicy` method), 86  
 forward () (`genrl.core.base.BaseValue` method), 86  
 forward () (`genrl.core.noise.NoisyLinear` method), 89  
 forward () (`genrl.core.policies.CNNPolicy` method),  
 forward () (`genrl.core.values.CnnCategoricalValue` method), 94  
 forward () (`genrl.core.values.CnnDuelingValue` method), 94  
 forward () (`genrl.core.values.CnnNoisyValue` method), 95  
 forward () (`genrl.core.values.CnnValue` method), 95  
 forward () (`genrl.core.values.MlpCategoricalValue` method), 96  
 forward () (`genrl.core.values.MlpDuelingValue` method), 96  
 FrameStack (class in `genrl.environments.frame_stack`), 79

framestack (*genrl.core.values.CnnCategoricalValue attribute*), 93  
 framestack (*genrl.core.values.CnnDuelingValue attribute*), 94

**G**

gamma (*genrl.agents.classical.qlearning.qlearning.QLearning attribute*), 56  
 gamma (*genrl.agents.classical.sarsa.sarsa.SARSA attribute*), 57  
 gamma (*genrl.agents.deep.a2c.a2c.A2C attribute*), 31  
 gamma (*genrl.agents.deep.ddpg.ddpg.DDPG attribute*), 33  
 gamma (*genrl.agents.deep.dqn.base.DQN attribute*), 35  
 gamma (*genrl.agents.deep.dqn.categorical.CategoricalDQN attribute*), 38  
 gamma (*genrl.agents.deep.dqn.double.DoubleDQN attribute*), 40  
 gamma (*genrl.agents.deep.dqn.dueling.DuelingDQN attribute*), 42  
 gamma (*genrl.agents.deep.dqn.noisy.NoisyDQN attribute*), 43  
 gamma (*genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN attribute*), 44  
 gamma (*genrl.agents.deep.ppo1.ppo1.PPO1 attribute*), 47  
 gamma (*genrl.agents.deep.sac.sac.SAC attribute*), 53  
 gamma (*genrl.agents.deep.td3.td3.TD3 attribute*), 51  
 GaussianMAB (*class in genrl.agents.bandits.multiarmed.gaussian\_mab*), 70  
 genrl.agents.bandits.contextual.base (*module*), 57  
 genrl.agents.bandits.contextual.bootstrap (*module*), 58  
 genrl.agents.bandits.contextual.common.base (*module*), 107  
 genrl.agents.bandits.contextual.common.baselines.actor\_critic (*module*), 108  
 genrl.agents.bandits.contextual.common.noise (*module*), 109  
 genrl.agents.bandits.contextual.common.transition\_policies (*module*), 110  
 genrl.agents.bandits.contextual.fixed (*module*), 59  
 genrl.agents.bandits.contextual.linpos (*module*), 60  
 genrl.agents.bandits.contextual.neural\_gated (*module*), 61  
 genrl.agents.bandits.contextual.neural\_lipschitz (*module*), 62  
 genrl.agents.bandits.contextual.neural\_noisy\_sampling (*module*), 63

genrl.agents.bandits.contextual.variational (module), 65  
 genrl.agents.bandits.multiarmed.base (module), 66  
 genrl.agents.bandits.multiarmed.bayesian (module), 67  
 genrl.agents.bandits.multiarmed.bernoulli\_mab (module), 68  
 genrl.agents.bandits.multiarmed.epsgreedy (module), 69  
 genrl.agents.bandits.multiarmed.gaussian\_mab (module), 70  
 genrl.agents.bandits.multiarmed.gradient (module), 70  
 genrl.agents.bandits.multiarmed.thompson (module), 71  
 genrl.agents.bandits.multiarmed.ucb (module), 72  
 genrl.agents.classical.qlearning.qlearning (module), 55  
 genrl.agents.classical.sarsa.sarsa (module), 56  
 genrl.agents.deep.a2c.a2c (module), 31  
 genrl.agents.deep.ddpg.ddpg (module), 33  
 genrl.agents.deep.dqn.base (module), 35  
 genrl.agents.deep.dqn.categorical (module), 38  
 genrl.agents.deep.dqn.double (module), 40  
 genrl.agents.deep.dqn.dueling (module), 41  
 genrl.agents.deep.dqn.noisy (module), 43  
 genrl.agents.deep.dqn.prioritized (module), 44  
 genrl.agents.deep.dqn.utils (module), 45  
 genrl.agents.deep.ppo1.ppo1 (module), 47  
 genrl.agents.deep.sac.sac (module), 53  
 genrl.agents.deep.td3.td3 (module), 51  
 genrl.agents.deep.vpg.vpg (module), 49  
 genrl.core.base (module), 85  
 genrl.core.buffers (module), 87  
 genrl.core.noise (module), 89  
 genrl.core.rollout\_storage (module), 91  
 genrl.core.values (module), 93  
 genrl.environments (module), 82  
 genrl.environments.action\_wrappers (module), 76  
 genrl.environments.atari\_preprocessing (module), 77  
 genrl.environments.atari\_wrappers (module), 78  
 genrl.environments.base\_wrapper (module), 78  
 genrl.environments.frame\_stack (module),

79  
genrl.environments.gym\_wrapper (*module*),  
80  
genrl.environments.suite (*module*), 80  
genrl.environments.time\_limit (*module*), 81  
genrl.environments.vec\_env (*module*), 76  
genrl.environments.vec\_env.monitor (*module*), 73  
genrl.environments.vec\_env.normalize  
(*module*), 74  
genrl.environments.vec\_env.utils (*module*), 74  
genrl.environments.vec\_env.vector\_envs  
(*module*), 74  
genrl.environments.vec\_env.wrappers  
(*module*), 76  
genrl.trainers.ClassicalTrainer (*module*),  
104  
genrl.trainers.DCBTrainer (*module*), 104  
genrl.trainers.MABTrainer (*module*), 105  
genrl.trainers.OffPolicyTrainer (*module*),  
102  
genrl.trainers.OnPolicyTrainer (*module*),  
101  
genrl.trainers.Trainer (*module*), 105  
genrl.utils.logger (*module*), 97  
genrl.utils.models (*module*), 100  
genrl.utils.utils (*module*), 98  
get () (genrl.core.rollout\_storage.RolloutBuffer  
*method*), 93  
get\_action () (genrl.agents.classical.qlearning.qlearning.QLearn  
*method*), 56  
get\_action () (genrl.agents.classical.sarsa.sarsa.SARSA  
*method*), 57  
get\_action () (genrl.core.actor\_critic.CNNActorCritic  
*method*), 83  
get\_action () (genrl.core.actor\_critic.MlpSingleActorMultiCritic  
*method*), 85  
get\_action () (genrl.core.base.BaseActorCritic  
*method*), 85  
get\_action () (genrl.core.base.BasePolicy  
*method*),  
86  
get\_actor\_critic\_from\_name () (in *module*  
genrl.core.actor\_critic), 85  
get\_alpha\_loss () (genrl.agents.deep.sac.sac.SAC  
*method*), 54  
get\_data () (genrl.agents.bandits.contextual.common.transition.TransitionDB  
*method*), 110  
get\_data\_for\_action ()  
(genrl.agents.bandits.contextual.common.transition.TransitionDB  
*method*), 110  
get\_env\_properties () (in *module*  
genrl.utils.utils), 99  
get\_greedy\_action ()  
(genrl.agents.deep.dqn.base.DQN  
*method*),  
36  
get\_greedy\_action ()  
(genrl.agents.deep.dqn.categorical.CategoricalDQN  
*method*), 39  
get\_hyperparams ()  
(genrl.agents.classical.qlearning.qlearning.QLearning  
*method*), 56  
get\_hyperparams ()  
(genrl.agents.deep.a2c.a2c.A2C  
*method*),  
32  
get\_hyperparams ()  
(genrl.agents.deep.ddpg.ddpg.DDPG  
*method*),  
34  
get\_hyperparams ()  
(genrl.agents.deep.dqn.base.DQN  
*method*),  
36  
get\_hyperparams ()  
(genrl.agents.deep.ppo1.ppo1.PPO1  
*method*),  
48  
get\_hyperparams () (genrl.agents.deep.sac.sac.SAC  
*method*), 54  
get\_hyperparams () (genrl.agents.deep.td3.td3.TD3  
*method*), 52  
get\_hyperparams ()  
(genrl.agents.deep.vpg.vpg.VPG  
*method*),  
50  
get\_log\_probs () (genrl.agents.deep.vpg.vpg.VPG  
*method*), 50  
get\_logger\_by\_name () (in *module*  
genrl.utils.logger), 98  
get\_logging\_params ()  
(genrl.agents.deep.a2c.a2c.A2C  
*method*),  
32  
get\_logging\_params ()  
(genrl.agents.deep.ddpg.ddpg.DDPG  
*method*),  
35  
get\_logging\_params ()  
(genrl.agents.deep.dqn.base.DQN  
*method*),  
36  
get\_logging\_params ()  
(genrl.agents.deep.ppo1.ppo1.PPO1  
*method*),  
49  
get\_logging\_params ()  
(genrl.agents.deep.sac.sac.SAC  
*method*),  
54  
get\_transitionDB\_params ()  
(genrl.agents.deep.td3.td3.TD3  
*method*),  
53  
get\_transitionDBj\_params ()  
(genrl.agents.deep.vpg.vpg.VPG  
*method*),  
50  
get\_model () (in *module* genrl.utils.utils), 99  
get\_model\_from\_name () (in *module*

```

        genrl.utils.models), 100
get_p_loss() (genrl.agents.deep.sac.sac.SAC
    method), 55
get_policy_from_name() (in module
    genrl.core.policies), 91
get_q_loss() (genrl.agents.deep.dqn.categorical.CategoricalDQN
    method), 39
get_q_loss() (genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN
    method), 45
get_q_values() (genrl.agents.deep.dqn.base.DQN
    method), 37
get_q_values() (genrl.agents.deep.dqn.categorical.CategoricalDQN
    method), 39
get_spaces() (genrl.environments.vec_env.vector_envs.SerialVecEnv
    method), 75
get_spaces() (genrl.environments.vec_env.vector_envs.SubProcessVecEnv
    method), 75
get_target_q_values() (genrl.agents.deep.dqn.base.DQN
    method), 37
get_target_q_values() (genrl.agents.deep.dqn.categorical.CategoricalDQN
    method), 40
get_target_q_values() (genrl.agents.deep.dqn.double.DoubleDQN
    method), 41
get_target_q_values() (genrl.agents.deep.sac.sac.SAC
    method), 55
get_traj_loss() (genrl.agents.deep.a2c.a2c.A2C
    method), 32
get_traj_loss() (genrl.agents.deep.ppo1.ppo1.PPO1
    method), 49
get_traj_loss() (genrl.agents.deep.vpg.vpg.VPG
    method), 50
get_value() (genrl.core.actor_critic.CNNActorCritic
    method), 83
get_value() (genrl.core.actor_critic.MlpSingleActorMultiCritic
    method), 85
get_value() (genrl.core.base.BaseActorCritic
    method), 86
get_value() (genrl.core.base.BaseValue method), 86
get_value_from_name() (in module
    genrl.core.values), 97
GradientMABAgent (class) in
    genrl.agents.bandits.multiarmed.gradient,
    70
GymEnv() (in module genrl.environments.suite), 81
GymWrapper (class) in
    genrl.environments.gym_wrapper), 80
H
hidden (genrl.core.actor_critic.MlpActorCritic at-
tribute), 83
hidden (genrl.core.actor_critic.MlpSingleActorMultiCritic
    attribute), 84
hidden (genrl.core.values.MlpDuelingValue attribute),
    96
HumanOutputFormat (class in genrl.utils.logger), 97
I
Integers (genrl.environments.vec_env.vector_envs.SerialVecEnv
    method), 75
in_features (genrl.core.noise.NoisyLinear attribute),
    89
IS_EMPTY (genrl.core.buffers.PrioritizedReplayBufferSamples
    attribute), 88
IS_EMPTY (genrl.utils.models.TabularModel
    method), 100
L
layers (genrl.agents.deep.a2c.a2c.A2C attribute), 31
layers (genrl.agents.deep.ddpg.ddpg.DDPG attribute),
    34
layers (genrl.agents.deep.dqn.categorical.CategoricalDQN
    attribute), 38
layers (genrl.agents.deep.dqn.double.DoubleDQN at-
tribute), 40
layers (genrl.agents.deep.dqn.dueling.DuelingDQN
    attribute), 42
layers (genrl.agents.deep.dqn.noisy.NoisyDQN at-
tribute), 43
layers (genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN
    attribute), 44
layers (genrl.agents.deep.ppo1.ppo1.PPO1 attribute),
    47
LazyFrames (class) in
    genrl.environments.frame_stack), 79
LinearPosteriorAgent (class) in
    genrl.agents.bandits.contextual.linpos), 60
load_model (in module
    genrl.trainers.OffPolicyTrainer), 103
load_model (in module
    genrl.trainers.OnPolicyTrainer), 102
load_model (in module genrl.trainers.Trainer), 106
load_weights() (genrl.agents.deep.a2c.a2c.A2C
    method), 33
load_weights() (genrl.agents.deep.dqn.base.DQN
    method), 37
load_weights() (genrl.agents.deep.ppo1.ppo1.PPO1
    method), 49
load_weights() (genrl.agents.deep.vpg.vpg.VPG
    method), 50
log_interval (in module
    genrl.trainers.OffPolicyTrainer), 103
log_interval (in module
    genrl.trainers.OnPolicyTrainer), 101
log_interval (in module genrl.trainers.Trainer), 105

```

log\_key (in module `genrl.trainers.OffPolicyTrainer`), 103  
log\_key (in module `genrl.trainers.OnPolicyTrainer`), 101  
log\_key (in module `genrl.trainers.Trainer`), 105  
log\_mode (in module `genrl.trainers.OffPolicyTrainer`), 103  
log\_mode (in module `genrl.trainers.OnPolicyTrainer`), 101  
log\_mode (in module `genrl.trainers.Trainer`), 105  
logdir (`genrl.utils.logger.Logger` attribute), 98  
logdir (in module `genrl.trainers.OffPolicyTrainer`), 103  
logdir (in module `genrl.trainers.OnPolicyTrainer`), 101  
logdir (in module `genrl.trainers.Trainer`), 105  
Logger (class in `genrl.utils.logger`), 97  
lr (`genrl.agents.classical.qlearning.QLearning` attribute), 56  
lr (`genrl.agents.classical.sarsa.sarsa.SARSA` attribute), 57  
lr\_policy (`genrl.agents.deep.a2c.a2c.A2C` attribute), 31  
lr\_policy (`genrl.agents.deep.ddpg.ddpg.DDPG` attribute), 34  
lr\_policy (`genrl.agents.deep.ppo1.ppo1.PPO1` attribute), 47  
lr\_policy (`genrl.agents.deep.sac.sac.SAC` attribute), 53  
lr\_policy (`genrl.agents.deep.td3.td3.TD3` attribute), 52  
lr\_value (`genrl.agents.deep.a2c.a2c.A2C` attribute), 31  
lr\_value (`genrl.agents.deep.ddpg.ddpg.DDPG` attribute), 34  
lr\_value (`genrl.agents.deep.dqn.base.DQN` attribute), 35  
lr\_value (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
lr\_value (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 40  
lr\_value (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 42  
lr\_value (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
lr\_value (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN` attribute), 44  
lr\_value (`genrl.agents.deep.ppo1.ppo1.PPO1` attribute), 48  
lr\_value (`genrl.agents.deep.sac.sac.SAC` attribute), 54  
lr\_value (`genrl.agents.deep.td3.td3.TD3` attribute), 52

**M**

MABAgent (class in `genrl.agents.bandits.multiarmed.base`), 66

max\_ep\_len (in module `genrl.trainers.OffPolicyTrainer`), 102  
max\_epsilon (`genrl.agents.deep.dqn.base.DQN` attribute), 36  
max\_epsilon (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
max\_epsilon (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 41  
max\_epsilon (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 42  
max\_epsilon (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
max\_epsilon (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN` attribute), 45  
max\_key\_len () (`genrl.utils.logger.HumanOutputFormat` method), 97  
max\_timesteps (in module `genrl.trainers.OffPolicyTrainer`), 103  
max\_timesteps (in module `genrl.trainers.OnPolicyTrainer`), 101  
max\_timesteps (in module `genrl.trainers.Trainer`), 106  
mean (`genrl.core.noise.ActionNoise` attribute), 89  
min\_epsilon (`genrl.agents.deep.dqn.base.DQN` attribute), 36  
min\_epsilon (`genrl.agents.deep.dqn.categorical.CategoricalDQN` attribute), 38  
min\_epsilon (`genrl.agents.deep.dqn.double.DoubleDQN` attribute), 41  
min\_epsilon (`genrl.agents.deep.dqn.dueling.DuelingDQN` attribute), 42  
min\_epsilon (`genrl.agents.deep.dqn.noisy.NoisyDQN` attribute), 43  
min\_epsilon (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN` attribute), 45  
mlp () (in module `genrl.utils.utils`), 99  
**DQNActorCritic** (class in `genrl.core.actor_critic`), 83  
**MlpCategoricalValue** (class in `genrl.core.values`), 95  
**MlpDuelingValue** (class in `genrl.core.values`), 96  
**MlpNoisyValue** (class in `genrl.core.values`), 96  
**MlpPolicy** (class in `genrl.core.policies`), 90  
**MlpSingleActorMultiCritic** (class in `genrl.core.actor_critic`), 84  
**PrioritizedReplayDQNE** (class in `genrl.core.values`), 96  
**Model** (class in `genrl.agents.bandits.contextual.common.base_model`), 107

**N**

n\_envs (`genrl.environments.vec_env.vector_envs.VecEnv` attribute), 76  
n\_envs (in module `genrl.trainers.OffPolicyTrainer`), 104  
n\_envs (in module `genrl.trainers.OnPolicyTrainer`), 102

n\_envs (in module `genrl.trainers.Trainer`), 106  
n\_episodes (`genrl.core.rollout_storage.RolloutReturn attribute`), 93  
network (`genrl.agents.deep.a2c.a2c.A2C attribute`), 31  
network (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 33  
network (`genrl.agents.deep.dqn.base.DQN attribute`), 35  
network (`genrl.agents.deep.dqn.categorical.CategoricalDQN attribute`), 38  
network (`genrl.agents.deep.dqn.double.DoubleDQN attribute`), 40  
network (`genrl.agents.deep.dqn.dueling.DuelingDQN attribute`), 41  
network (`genrl.agents.deep.dqn.noisy.NoisyDQN attribute`), 43  
network (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN attribute`), 44  
network (`genrl.agents.deep.ppo1.ppo1.PPO1 attribute`), 47  
network (`genrl.agents.deep.sac.sac.SAC attribute`), 53  
network (`genrl.agents.deep.td3.td3.TD3 attribute`), 51  
NeuralBanditModel (class in `genrl.agents.bandits.contextual.common.neural`), 109  
NeuralGreedyAgent (class in `genrl.agents.bandits.contextual.neural_greedy`), 61  
NeuralLinearPosteriorAgent (class in `genrl.agents.bandits.contextual.neural_lipos`), 62  
NeuralNoiseSamplingAgent (class in `genrl.agents.bandits.contextual.neural_noise_sampling`), 63  
next\_observations (`genrl.core.rollout_storage.ReplayBufferSamples attribute`), 92  
next\_states (`genrl.core.buffers.PrioritizedReplayBufferSamples attribute`), 88  
next\_states (`genrl.core.buffers.ReplayBufferSamples attribute`), 89  
noise (`genrl.agents.deep.a2c.a2c.A2C attribute`), 31  
noise (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 34  
noise (`genrl.agents.deep.td3.td3.TD3 attribute`), 52  
noise\_std (`genrl.agents.deep.a2c.a2c.A2C attribute`), 31  
noise\_std (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 34  
noise\_std (`genrl.agents.deep.td3.td3.TD3 attribute`), 52  
noisy\_layers (`genrl.agents.deep.dqn.categorical.CategoricalDQN attribute`), 39  
noisy\_layers (`genrl.agents.deep.dqn.noisy.NoisyDQN attribute`), 43  
noisy\_layers (`genrl.core.values.CnnCategoricalValue attribute`), 94  
noisy\_layers (`genrl.core.values.CnnNoisyValue attribute`), 95  
noisy\_layers (`genrl.core.values.MlpCategoricalValue attribute`), 95  
noisy\_mlp () (in module `genrl.utils.utils`), 99  
NoisyDQN (class in `genrl.agents.deep.dqn.noisy`), 43  
NoisyLinear (class in `genrl.core.noise`), 89  
NoopReset (class in `genrl.environments.atari_wrappers`), 78  
NormalActionNoise (class in `genrl.core.noise`), 90  
num\_atoms (`genrl.agents.deep.dqn.categorical.CategoricalDQN attribute`), 39  
num\_atoms (`genrl.core.values.CnnCategoricalValue attribute`), 94  
num\_atoms (`genrl.core.values.CnnNoisyValue attribute`), 95  
num\_atoms (`genrl.core.values.MlpCategoricalValue attribute`), 96  
num\_critics (`genrl.core.actor_critic.MlpSingleActorMultiCritic attribute`), 84  
O  
obs\_shape (`genrl.environments.gym_wrapper.GymWrapper attribute`), 80  
obs\_shape (`genrl.environments.vec_env.vector_envs.VecEnv attribute`), 76  
observation\_spaces (`genrl.environments.vec_env.vector_envs.VecEnv attribute`), 76  
Observations (`genrl.core.rollout_storage.ReplayBufferSamples attribute`), 92  
observations (`genrl.core.rollout_storage.RolloutBufferSamples attribute`), 93  
off\_policy (in module `genrl.trainers.OffPolicyTrainer`), 103  
off\_policy (in module `genrl.trainers.OnPolicyTrainer`), 101  
off\_policy (in module `genrl.trainers.Trainer`), 106  
old\_log\_prob (`genrl.core.rollout_storage.RolloutBufferSamples attribute`), 93  
old\_values (`genrl.core.rollout_storage.RolloutBufferSamples attribute`), 93  
OrnsteinUhlenbeckActionNoise (class in `genrl.core.noise`), 90  
out\_features (`genrl.core.noise.NoisyLinear attribute`), 89  
P  
policy\_frequency (`genrl.agents.deep.td3.td3.TD3 attribute`), 52

policy\_layers (`genrl.agents.deep.sac.sac.SAC attribute`), 53  
 policy\_layers (`genrl.agents.deep.td3.td3.TD3 attribute`), 51  
 polyak (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 34  
 polyak (`genrl.agents.deep.sac.sac.SAC attribute`), 54  
 polyak (`genrl.agents.deep.td3.td3.TD3 attribute`), 52  
 pos (`genrl.core.buffers.PrioritizedBuffer attribute`), 87  
 PPO1 (`class in genrl.agents.deep.ppo1.ppo1`), 47  
 prioritized\_q\_loss () (in module `genrl.agents.deep.dqn.utils`), 47  
 PrioritizedBuffer (`class in genrl.core.buffers`), 87  
 PrioritizedReplayBufferSamples (`class in genrl.core.buffers`), 87  
 PrioritizedReplayDQN (`class in genrl.agents.deep.dqn.prioritized`), 44  
 probability\_hist (`genrl.agents.bandits.multiarmed.attribute`), 70  
 push () (`genrl.core.buffers.PrioritizedBuffer method`), 87  
 push () (`genrl.core.buffers.PushReplayBuffer method`), 88  
 push () (`genrl.core.buffers.ReplayBuffer method`), 88  
 PushReplayBuffer (`class in genrl.core.buffers`), 88

**Q**

QLearning (class in `genrl.agents.classical.qlearning.qlearning`), 55  
 quality (`genrl.agents.bandits.multiarmed.bayesian.BayesianUCBMABAgent attribute`), 68  
 quality (`genrl.agents.bandits.multiarmed.epsgreedy.EpsGreedyMABAgent attribute`), 69  
 quality (`genrl.agents.bandits.multiarmed.gradient.GradientMABAgent attribute`), 70  
 quality (`genrl.agents.bandits.multiarmed.thompson.ThompsonSamplingMABAgent attribute`), 71  
 quality (`genrl.agents.bandits.multiarmed.ucb.UCBMABAgent attribute`), 72

**R**

regret (`genrl.agents.bandits.multiarmed.base.MABAgent attribute`), 66  
 regret\_hist (`genrl.agents.bandits.multiarmed.base.MABAgent attribute`), 66  
 render (`genrl.agents.deep.a2c.a2c.A2C attribute`), 32  
 render (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 34  
 render (`genrl.agents.deep.dqn.base.DQN attribute`), 36  
 render (`genrl.agents.deep.dqn.categorical.CategoricalDQN attribute`), 39  
 render (`genrl.agents.deep.dqn.double.DoubleDQN attribute`), 41

render (`genrl.agents.deep.dqn.dueling.DuelingDQN attribute`), 42  
 render (`genrl.agents.deep.dqn.noisy.NoisyDQN attribute`), 44  
 render (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN attribute`), 45  
 render (`genrl.agents.deep.ppo1.ppo1.PPO1 attribute`), 48  
 render (`genrl.agents.deep.sac.sac.SAC attribute`), 54  
 render (`genrl.agents.deep.td3.td3.TD3 attribute`), 52  
 render (in module `genrl.trainers.OffPolicyTrainer`), 103  
 render (in module `genrl.trainers.OnPolicyTrainer`), 102  
 render (in module `genrl.trainers.Trainer`), 106  
 render () (`genrl.environments.base_wrapper.BaseWrapper method`), 78  
 render () (`genrl.environments.gym_wrapper.GymWrapper gradient.MABAgent`), 75  
 render () (`genrl.environments.vec_env.vector_envs.SerialVecEnv method`), 75  
 replay\_size (`genrl.agents.deep.ddpg.ddpg.DDPG attribute`), 34  
 replay\_size (`genrl.agents.deep.dqn.base.DQN attribute`), 35  
 replay\_size (`genrl.agents.deep.dqn.categorical.CategoricalDQN attribute`), 38  
 replay\_size (`genrl.agents.deep.dqn.double.DoubleDQN attribute`), 40  
 replay\_size (`genrl.agents.deep.dqn.dueling.DuelingDQN attribute`), 42  
 replay\_size (`genrl.agents.deep.dqn.noisy.NoisyDQN attribute`), 43  
 replay\_size (`genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN attribute`), 44  
 replay\_size (`genrl.agents.deep.sac.sac.SAC attribute`), 54  
 replay\_size (`genrl.agents.deep.td3.td3.TD3 attribute`), 52  
 ReplayBuffer (`class in genrl.core.buffers`), 88  
 ReplayBufferSamples (`class in genrl.core.buffers`), 88  
 ReplayBufferSamples (class in `genrl.core.rollout_storage`), 92  
 RescaleAction (class in `genrl.environments.action_wrappers`), 76  
 reset () (genrl.core.noise.NormalActionNoise method), 90  
 reset () (genrl.core.noise.OrnsteinUhlenbeckActionNoise method), 90  
 reset () (genrl.core.rollout\_storage.BaseBuffer method), 91  
 reset () (genrl.core.rollout\_storage.RolloutBuffer

```

        method), 93
reset () (genrl.environments.atari_preprocessing.AtariPreprocessing RolloutBufferSamples (class in
        method), 77 genrl.core.rollout_storage), 93
reset () (genrl.environments.atari_wrappers.FireReset RolloutReturn (class in genrl.core.rollout_storage),
        method), 78 93
reset () (genrl.environments.atari_wrappers.NoopReset round () (genrl.utils.logger.HumanOutputFormat
        method), 78 method), 97
reset () (genrl.environments.base_wrapper.BaseWrapper run_num (in module genrl.trainers.OffPolicyTrainer),
        method), 78 103
reset () (genrl.environments.frame_stack.FrameStack run_num (in module genrl.trainers.OnPolicyTrainer),
        method), 79 102
reset () (genrl.environments.gym_wrapper.GymWrapper run_num (in module genrl.trainers.Trainer), 106
        method), 80 RunningMeanStd (class in
        method), 80 genrl.environments.vec_env.utils), 74
reset () (genrl.environments.time_limit.AtariTimeLimit SAC (class in genrl.agents.deep.sac.sac), 53
        method), 81 sac (genrl.core.actor_critic.MlpActorCritic attribute),
        method), 82 83
reset () (genrl.environments.vec_env.monitor.VecMonitor sac (genrl.core.actor_critic.MlpSingleActorMultiCritic
        method), 73 attribute), 84
reset () (genrl.environments.vec_env.normalize.VecNormalize safe_mean () (in module genrl.utils.utils), 100
        method), 74 sample () (genrl.core.buffers.PrioritizedBuffer
        method), 75 method), 87
reset () (genrl.environments.vec_env.vector_envs.SerialVecEnv SampleVecEnv (genrl.core.buffers.PushReplayBuffer
        method), 75 method), 88
reset () (genrl.environments.vec_env.vector_envs.SubProcessVecEnv sample () (genrl.core.buffers.ReplayBuffer
        method), 75 method), 88
reset () (genrl.environments.vec_env.vector_envs.VecEnv sample () (genrl.core.buffers.ReplayBuffer
        method), 76 method), 88
reset () (genrl.environments.vec_env.wrappers.VecEnvWrapper sample () (genrl.core.rollout_storage.BaseBuffer
        method), 91 method), 91
reset_noise () (genrl.core.noise.NoisyLinear sample () (genrl.environments.gym_wrapper.GymWrapper
        method), 89 method), 80
reset_noise () (genrl.core.values.MlpNoisyValue sample () (genrl.environments.vec_env.vector_envs.VecEnv
        method), 96 method), 76
reset_parameters () sample () (genrl.utils.models.TabularModel
        (genrl.agents.bandits.contextual.common.bayesian.BayesianLinear attribute), 56
        method), 108 save_interval (in module
        method), 90 genrl.trainers.OffPolicyTrainer), 103
returns (genrl.core.rollout_storage.RolloutBufferSamples save_interval (in module
        attribute), 93 genrl.trainers.OnPolicyTrainer), 101
        save_interval (in module genrl.trainers.Trainer),
reward_hist (genrl.agents.bandits.multiarmed.base.MABAgent 106
        attribute), 66 save_model (in module
        attribute), 88 genrl.trainers.OffPolicyTrainer), 103
        save_model (in module
rewards (genrl.core.buffers.PrioritizedReplayBufferSamples genrl.trainers.OnPolicyTrainer), 102
        attribute), 88 save_model (in module genrl.trainers.Trainer), 106
        attribute), 89 seed (genrl.agents.deep.a2c.a2c.A2C attribute), 32
        attribute), 92 seed (genrl.agents.deep.ddpg.ddpg.DDPG attribute), 34
rollout_size (genrl.agents.deep.a2c.a2c.A2C save_model (in module genrl.agents.deep.dqn.base.DQN
        attribute), 31 attribute), 36
        attribute), 48 seed (genrl.agents.deep.dqn.categorical.CategoricalDQN
        attribute), 39
RolloutBuffer (class in genrl.core.rollout_storage), seed (genrl.agents.deep.dqn.double.DoubleDQN
        92 attribute), 41

```

```
seed (genrl.agents.deep.dqn.dueling.DuelingDQN at- select_action () (genrl.agents.deep.sac.sac.SAC
tribute), 42 method), 55
seed (genrl.agents.deep.dqn.noisy.NoisyDQN attribute), select_action () (genrl.agents.deep.vpg.vpg.VPG
44 method), 50
seed (genrl.agents.deep.dqn.prioritized.PrioritizedReplayDQN nialVecEnv (class
attribute), 45 in
genrl.environments.vec_env.vector_envs),
74
seed (genrl.agents.deep.ppo1.ppo1.PPO1 attribute), 48 set_seeds () (in module genrl.utils.utils), 100
seed (genrl.agents.deep.sac.sac.SAC attribute), 54 shape (genrl.environments.frame_stack.LazyFrames at-
seed (genrl.agents.deep.td3.td3.TD3 attribute), 52 tribute), 79
seed (in module genrl.trainers.OffPolicyTrainer), 104 size () (genrl.core.rollout_storage.BaseBuffer
seed (in module genrl.trainers.OnPolicyTrainer), 102 method), 91
seed (in module genrl.trainers.Trainer), 106 start_update (in module
seed () (genrl.environments.base_wrapper.BaseWrapper genrl.trainers.OffPolicyTrainer), 102
method), 79 state_dim (genrl.core.actor_critic.MlpActorCritic at-
seed () (genrl.environments.gym_wrapper.GymWrapper method), 80 tribute), 83
seed () (genrl.environments.vec_env.vector_envs.SubProcessVecEnv dim (genrl.core.actor_critic.MlpSingleActorMultiCritic
method), 75 attribute), 84
seed () (genrl.environments.vec_env.vector_envs.VecEnv state_dim (genrl.core.values.CnnNoisyValue at-
method), 76 tribute), 94
select_action () (genrl.agents.bandits.contextual.bases.DCBAgent m (genrl.core.values.MlpCategoricalValue at-
attribute), 57 tribute), 95
select_action () (genrl.agents.bandits.contextual.bootstrap_mixedAgents.BootstrapMixedAgents.MlpDuelingValue
method), 59 attribute), 96
select_action () (genrl.agents.bandits.contextual.fixed_FixedAgent genrl.core.buffers.PrioritizedReplayBufferSamples
method), 59 attribute), 88
select_action () (genrl.agents.bandits.contextual.linpost_LinPostLinearPosteridgAgentore.buffers.ReplayBufferSamples
method), 60 attribute), 89
select_action () (genrl.agents.bandits.contextual.neural_gatedNoise.NeuralGatedNoise attribute), 89
method), 61 std_init (genrl.core.noise.Noisy attribute), 89
select_action () (genrl.agents.bandits.contextual.neural_gatedNoise.NeuralGatedNoise attribute), 89
method), 63 m (genrl.core.buffers.PrioritizedReplayBufferSamples
method), 77
select_action () (genrl.agents.bandits.contextual.neural_noopReset_NoopReset attribute), 89
method), 64 std_init (genrl.core.buffers.PrioritizedReplayBufferSamples
method), 78
select_action () (genrl.agents.bandits.contextual.variational_VariationalAgent genrl.environments.base_wrapper.BaseWrapper
method), 65 attribute), 79
select_action () (genrl.agents.bandits.multiarmed.base_MARAgent genrl.environments.frame_stack.FrameStack
method), 66 attribute), 79
select_action () (genrl.agents.bandits.multiarmed.bayesian_BayesianMABAgent gym_wrapper.GymWrapper
method), 68 attribute), 80
select_action () (genrl.agents.bandits.multiarmed.epsGreedy_EpsGreedyMABAgent time_limit.AtariTimeLimit
method), 69 attribute), 81
select_action () (genrl.agents.bandits.multiarmed.gradient_GradientMABAgent environments.time_limit.TimeLimit
method), 70 attribute), 82
select_action () (genrl.agents.bandits.multiarmed.thompsonSampling_ThompsonSamplingMABAgent monitor.VecMonitor
method), 71 attribute), 73
select_action () (genrl.agents.bandits.multiarmed.ucb_UCBMAgent environments.vec_env.normalize.VecNormalize
method), 72 attribute), 74
select_action () (genrl.agents.deep.a2c.a2c.A2C step () (genrl.environments.vec_env.vector_envs.SerialVecEnv
method), 33 attribute), 75
select_action () (genrl.agents.deep.dqn.base.DQN step () (genrl.environments.vec_env.vector_envs.SubProcessVecEnv
method), 37 attribute), 75
select_action () (genrl.agents.deep.ppo1.ppo1.PPO1 step () (genrl.environments.vec_env.vector_envs.VecEnv
method), 49 attribute), 76
```

step () (genrl.environments.vec\_env.wrappers.VecEnvWrapper method), 58  
 step () (genrl.utils.models.TabularModel method), 100  
 SubProcessVecEnv (class in genrl.environments.vec\_env.vector\_envs), 75  
 swap\_and\_flatten () (genrl.core.rollout\_storage.BaseBuffer static method), 91

**T**

TabularModel (class in genrl.utils.models), 100  
 TD3 (class in genrl.agents.deep.td3.td3), 51  
 temp (genrl.agents.bandits.multiarmed.gradient.GradientMABAgent attribute), 71  
 TensorboardLogger (class in genrl.utils.logger), 98  
 ThompsonSamplingMABAgent (class in genrl.agents.bandits.multiarmed.thompson), 71  
 TimeLimit (class in genrl.environments.time\_limit), 82  
 to\_torch () (genrl.core.rollout\_storage.BaseBuffer method), 91  
 train\_model () (genrl.agents.bandits.contextual.common.PasteModelModel method), 107  
 TransitionDB (class in genrl.agents.bandits.contextual.common.transition), 110

**U**

UCBMABAgent (class in genrl.agents.bandits.multiarmed.ucb), 72  
 update () (genrl.agents.classical.qlearning.QLearning method), 35  
 update () (genrl.agents.classical.sarsa.sarsa.SARSA method), 57  
 update () (genrl.environments.vec\_env.utils.RunningMeanStd method), 74  
 update\_db () (genrl.agents.bandits.contextual.bootstrap\_neural.BootstrapNeuralAgent method), 59  
 update\_db () (genrl.agents.bandits.contextual.fixed.FixedAgent method), 59  
 update\_db () (genrl.agents.bandits.contextual.linpos.LinearPosterioAgent method), 51  
 update\_db () (genrl.agents.bandits.contextual.neural\_greedy.NeuralGreedyAgent method), 61  
 update\_db () (genrl.agents.bandits.contextual.neural\_linpos.NeuralLinposAgent method), 63  
 update\_db () (genrl.agents.bandits.contextual.noise\_sampling.NeuralNoiseSamplingAgent method), 64  
 update\_db () (genrl.agents.bandits.contextual.variational.VariationalAgent method), 65  
 update\_interval (in module genrl.trainers.OffPolicyTrainer), 103  
 update\_parameters () (genrl.agents.bandits.contextual.base.DCBAgent

method), 58  
 update\_params () (genrl.agents.bandits.contextual.bootstrap\_neural.BootstrapNeuralAgent method), 59  
 update\_params () (genrl.agents.bandits.contextual.fixed.FixedAgent method), 59  
 update\_params () (genrl.agents.bandits.contextual.linpos.LinearPosterioAgent method), 60  
 update\_params () (genrl.agents.bandits.contextual.neural\_greedy.NeuralGreedyAgent method), 61  
 update\_params () (genrl.agents.bandits.contextual.neural\_linpos.NeuralLinposAgent method), 63  
 update\_params () (genrl.agents.bandits.contextual.noise\_sampling.NeuralNoiseSamplingAgent method), 64  
 update\_params () (genrl.agents.bandits.contextual.variational.VariationalAgent method), 65  
 update\_params () (genrl.agents.bandits.contextual.variational.VariationalAgent method), 65  
 update\_params () (genrl.agents.bandits.multiarmed.base.MABAgent method), 67  
 update\_params () (genrl.agents.bandits.multiarmed.bayesian.BayesianAgent method), 68  
 update\_params () (genrl.agents.bandits.multiarmed.epsgreedy.EpsGreedyAgent method), 69  
 update\_params () (genrl.agents.bandits.multiarmed.gradient.GradientMABAgent method), 71  
 update\_params () (genrl.agents.bandits.multiarmed.thompson.ThompsonSamplingMABAgent method), 72  
 update\_params () (genrl.agents.bandits.multiarmed.ucb.UBCMABAgent method), 72  
 update\_params () (genrl.agents.deep.a2c.a2c.A2CAgent method), 33  
 update\_params () (genrl.agents.deep.ddpg.ddpg.DDPG method), 35  
 update\_params () (genrl.agents.deep.dqn.base.DQN method), 37  
 update\_params () (genrl.agents.deep.ppo1.ppo1.PPO1 method), 49  
 update\_params () (genrl.agents.deep.sac.sac.SAC method), 53  
 update\_params () (genrl.agents.deep.td3.td3.TD3 method), 55  
 update\_params () (genrl.agents.deep.vpg.vpg.VPG method), 51  
 update\_params\_before\_select\_action ()  
 update\_priorities ()  
 update\_prioritized\_buffer (genrl.core.buffer.PrioritizedBuffer method), 87  
 update\_target\_node (genrl.core.buffer.PrioritizedBuffer method), 87  
 update\_target\_model ()  
 (genrl.agents.deep.sac.sac.SAC method), 55  
 use\_dropout (genrl.agents.bandits.contextual.common.base\_model.Model attribute), 107

use\_dropout (*genrl.agents.bandits.contextual.common.bayesian.BayesianNNBanditModel attribute*), 109  
use\_dropout (*genrl.agents.bandits.contextual.common.neural.NeuralBanditModel attribute*), 109

**V**

v\_max (*genrl.agents.deep.dqn.categorical.CategoricalDQN attribute*), 39  
v\_min (*genrl.agents.deep.dqn.categorical.CategoricalDQN attribute*), 39  
val\_type (*genrl.core.actor\_critic.MlpActorCritic attribute*), 83  
val\_type (*genrl.core.actor\_critic.MlpSingleActorMultiCritic attribute*), 84  
value\_coeff (*genrl.agents.deep.a2c.a2c.A2C attribute*), 32  
value\_coeff (*genrl.agents.deep.ppo1.ppo1.PPO1 attribute*), 48  
value\_layers (*genrl.agents.deep.dqn.base.DQN attribute*), 35  
value\_layers (*genrl.agents.deep.sac.sac.SAC attribute*), 53  
value\_layers (*genrl.agents.deep.td3.td3.TD3 attribute*), 52  
VariationalAgent (class in *genrl.agents.bandits.contextual.variational*), 65  
VecEnv (class in *genrl.environments.vec\_env.vector\_envs*), 75  
VecEnvWrapper (class in *genrl.environments.vec\_env.wrappers*), 76  
VecMonitor (class in *genrl.environments.vec\_env.monitor*), 73  
VecNormalize (class in *genrl.environments.vec\_env.normalize*), 74  
VectorEnv () (in module *genrl.environments.suite*), 81  
VPG (class in *genrl.agents.deep.vpg.vpg*), 49

**W**

warmup\_steps (in module *genrl.trainers.OffPolicyTrainer*), 102  
weights (*genrl.core.buffers.PrioritizedReplayBufferSamples attribute*), 88  
worker () (in module *genrl.environments.vec\_env.vector\_envs*), 76  
write () (*genrl.utils.logger.CSVLogger method*), 97  
write () (*genrl.utils.logger.HumanOutputFormat method*), 97  
write () (*genrl.utils.logger.Logger method*), 98  
write () (*genrl.utils.logger.TensorboardLogger method*), 98  
write\_to\_file () (*genrl.utils.logger.HumanOutputFormat method*), 97