
RLzoo

Release 1.0.3

May 21, 2020

1	Installation	3
2	Quick Start	5
3	Configurations Overview	9
4	API	11
5	DQN and Variants	13
6	VPG	15
7	AC	19
8	A3C	21
9	DDPG	23
10	TD3	27
11	SAC	31
12	TRPO	35
13	PPO	39
14	DPPO	45
15	Basic Networks	51
16	Policy Networks	53
17	Value Networks	57
18	Replay Buffer	61
19	Distributions	65
20	Environment Wrappers	67

21 Environment List	71
22 Math Utilities	89
23 Common Utilities	91
24 DRL Book	93
25 DRL Tutorial	95
26 Contributing	97
27 Citation	99
Python Module Index	101
Index	103



RLzoo is a collection of the most practical reinforcement learning algorithms, frameworks and applications, released on [Github](#) in November 2019. It is implemented with Tensorflow 2.0 and API of neural network layers in Tensor-Layer 2, to provide a hands-on fast-developing approach for reinforcement learning practices and benchmarks. It supports basic toy-test environments like [OpenAI Gym](#) and [DeepMind Control Suite](#) with very simple configurations. Moreover, RLzoo supports robot learning benchmark environment [RLBench](#) based on Vrep/Pyrep simulator. Other large-scale distributed training framework for more realistic scenarios with Unity 3D, Mujoco, Bullet Physics, etc, will be supported in the future.

We also provide novices friendly [DRL Tutorials](#) for algorithms implementation, where each algorithm is implemented in an individual script. The tutorials serve as code examples for our Springer textbook [Deep Reinforcement Learning: Fundamentals, Research and Applications](#) , you can get the free PDF if your institute has Springer license.

CHAPTER 1

Installation

RLzoo generally requires Python \geq 3.5. Also if you want to use DeepMind Control Suite environment, Python 3.6 will be required.

Direct installation:

```
1 pip3 install rlzoo --upgrade
```

Install from the source code on github:

```
1 git clone https://github.com/tensorlayer/RLzoo.git
2 cd RLzoo
3 pip3 install .
```


2.1 Simple Usage

Open `./run_rlzoo.py`:

```
1 from rlzoo.common.env_wrappers import build_env
2 from rlzoo.common.utils import call_default_params
3 from rlzoo.algorithms import TD3
4 # choose an algorithm
5 AlgName = 'TD3'
6 # select a corresponding environment type
7 EnvType = 'classic_control'
8 # chose an environment
9 EnvName = 'Pendulum-v0'
10 # build an environment with wrappers
11 env = build_env(EnvName, EnvType)
12 # call default parameters for the algorithm and learning process
13 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
14 # instantiate the algorithm
15 alg = eval(AlgName+'(**alg_params)')
16 # start the training
17 alg.learn(env=env, mode='train', render=False, **learn_params)
18 # test after training
19 alg.learn(env=env, mode='test', render=True, **learn_params)
```

Run the example:

```
python run_rlzoo.py
```

Choices for AlgName: 'DQN', 'AC', 'A3C', 'DDPG', 'TD3', 'SAC', 'PG', 'TRPO', 'PPO', 'DPPO'

Choices for EnvType: 'atari', 'box2d', 'classic_control', 'mujoco', 'robotics', 'dm_control', 'rlbench'

Choices for EnvName refers to *List of Supported Environments in RLzoo*

2.2 Another Usage

For providing more flexibility, we provide another usage example of RLzoo with more explicit configurations as follows, where the users can pass in customized networks and optimizers, etc.

```

1  import gym
2  from rlzoo.common.utils import make_env, set_seed
3  from rlzoo.algorithms import AC
4  from rlzoo.common.value_networks import ValueNetwork
5  from rlzoo.common.policy_networks import StochasticPolicyNetwork
6
7  ''' load environment '''
8  env = gym.make('CartPole-v0').unwrapped
9  obs_space = env.observation_space
10 act_space = env.action_space
11 # reproducible
12 seed = 2
13 set_seed(seed, env)
14
15 ''' build networks for the algorithm '''
16 num_hidden_layer = 4 #number of hidden layers for the networks
17 hidden_dim = 64 # dimension of hidden layers for the networks
18 with tf.name_scope('AC'):
19     with tf.name_scope('Critic'):
20         # choose the critic network, can be replaced with customized network
21         critic = ValueNetwork(obs_space, hidden_dim_list=num_hidden_layer * _
↳[hidden_dim])
22     with tf.name_scope('Actor'):
23         # choose the actor network, can be replaced with customized network
24         actor = StochasticPolicyNetwork(obs_space, act_space, hidden_dim_
↳list=num_hidden_layer * [hidden_dim], output_activation=tf.nn.tanh)
25     net_list = [actor, critic] # list of the networks
26
27 ''' choose optimizers '''
28 a_lr, c_lr = 1e-4, 1e-2 # a_lr: learning rate of the actor; c_lr: learning rate of _
↳the critic
29 a_optimizer = tf.optimizers.Adam(a_lr)
30 c_optimizer = tf.optimizers.Adam(c_lr)
31 optimizers_list=[a_optimizer, c_optimizer] # list of optimizers
32
33 # initialize the algorithm model, with algorithm parameters passed in
34 model = AC(net_list, optimizers_list)
35 '''
36 full list of arguments for the algorithm
37 -----
38 net_list: a list of networks (value and policy) used in the algorithm, from common_
↳functions or customization
39 optimizers_list: a list of optimizers for all networks and differentiable variables
40 gamma: discounted factor of reward
41 action_range: scale of action values
42 '''
43
44 # start the training process, with learning parameters passed in
45 model.learn(env, train_episodes=500, max_steps=200,
46             save_interval=50, mode='train', render=False)
47 '''
48 full list of parameters for training

```

(continues on next page)

(continued from previous page)

```
49 -----
50 env: learning environment
51 train_episodes: total number of episodes for training
52 test_episodes: total number of episodes for testing
53 max_steps: maximum number of steps for one episode
54 save_interval: time steps for saving the weights and plotting the results
55 mode: 'train' or 'test'
56 render: if true, visualize the environment
57 '''
58
59 # test after training
60 model.learn(env, test_episodes=100, max_steps=200, mode='test', render=True)
```

2.3 Interactive Configurations

We also provide an interactive learning configuration with Jupyter Notebook and *ipywidgets*, where you can select the algorithm, environment, and general learning settings with simple clicking on dropdown lists and sliders! A video demonstrating the usage is as following. The interactive mode can be used with [rlzoo/interactive/main.ipynb](#) by running `$ jupyter notebook` to open it.

3.1 Supported DRL Algorithms

Generally RLzoo supports following DRL algorithms:

Value-based methods

- Deep Q-Networks (DQN)
- Double DQN
- Dueling DQN
- Prioritized Experience Replay (PER)
- Retrace
- Noisy DQN
- Distributed DQN

Policy-based methods

- Vanilla Policy Gradient (VPG)
- Trust Region Policy Optimization (TRPO)
- Proximal Policy Optimization (PPO)
- Distributed PPO (DPPO)

Actor-critic methods

- Actor-Critic (AC)
- Asynchronous Advantage Actor-Critic (A3C)
- Deep Deterministic Policy Gradient (DDPG)
- Twin Delayed DDPG (TD3)
- Soft Actor-Critic (SAC)

3.2 Supported Environments

Generally RLzoo supports following environments for DRL:

- **OpenAI Gym**
 - Atari
 - Box2D
 - Classic Control
 - MuJoCo
 - Robotics
- **DeepMind Control Suite**
- **RLBench**

Full list of specific names of environments supported in RLzoo can be checked in [List of Supported Environments in RLzoo](#).

3.3 Supported Configurations

Not all configurations (specific RL algorithm on specific environment) are supported in RLzoo, as in other libraries. The supported configurations for RL algorithms with corresponding environments in RLzoo are listed in the following table.

Algorithms	Action Space	Policy	Update	Envs
DQN (double, dueling, PER)	Discrete Only	NA	Off-policy	Atari, Classic Control
AC	Discrete/Continuous	Stochastic	On-policy	All
PG	Discrete/Continuous	Stochastic	On-policy	All
DDPG	Continuous	Deterministic	Off-policy	Classic Control, Box2D, MuJoCo, Robotics, DeepMind Control, RLBench
TD3	Continuous	Deterministic	Off-policy	Classic Control, Box2D, MuJoCo, Robotics, DeepMind Control, RLBench
SAC	Continuous	Stochastic	Off-policy	Classic Control, Box2D, MuJoCo, Robotics, DeepMind Control, RLBench
A3C	Discrete/Continuous	Stochastic	On-policy	Atari, Classic Control, Box2D, MuJoCo, Robotics, DeepMind Control
PPO	Discrete/Continuous	Stochastic	On-policy	All
DPPO	Discrete/Continuous	Stochastic	On-policy	Atari, Classic Control, Box2D, MuJoCo, Robotics, DeepMind Control
TRPO	Discrete/Continuous	Stochastic	On-policy	All

4.1 make_env()

It can be used as:

```
1 env = build_env(EnvName, EnvType)
```

4.2 call_default_params()

It can be used as:

```
1 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
```

The `call_default_params` returns the hyper-parameters stored in two dictionaries `alg_params` and `learn_params`, which can be printed to see what are contained inside. Hyper-parameters in these two dictionaries can also be changed by users before instantiating the agent and starting the learning process.

If you want to know exactly where the default hyper-parameters come from, they are stored in an individual Python script as `default.py` in each algorithm file in `./rlzoo/algorithms/`.

4.3 alg.learn()

It can be used as:

```
1 # start the training
2 alg.learn(env=env, mode='train', render=False, **learn_params)
3 # test after training
4 alg.learn(env=env, mode='test', render=True, **learn_params)
```

where the `alg` is an instantiation of DRL algorithm in RLzoo.

5.1 Example

```

1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import DQN
4
5  AlgName = 'DQN'
6  EnvName = 'PongNoFrameskip-v4'
7  EnvType = 'atari'
8
9  # EnvName = 'CartPole-v1'
10 # EnvType = 'classic_control' # the name of env needs to match the type of env
11
12 env = build_env(EnvName, EnvType)
13 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
14 alg = eval(AlgName+'(**alg_params)')
15 alg.learn(env=env, mode='train', **learn_params)
16 alg.learn(env=env, mode='test', render=True, **learn_params)

```

5.2 Deep Q-Networks

class rlzoo.algorithms.dqn.dqn.**DQN**(*net_list, optimizers_list, double_q, dueling, buffer_size, prioritized_replay, prioritized_alpha, prioritized_beta0*)

Papers:

Mnih V, Kavukcuoglu K, Silver D, et al. Human-level control through deep reinforcement learning[J]. Nature, 2015, 518(7540): 529.

Hessel M, Modayil J, Van Hasselt H, et al. Rainbow: Combining Improvements in Deep Reinforcement Learning[J]. 2017.

get_action (*obv, eps=0.2*)

get_action_greedy (*obv*)

learn (*env*, *mode*='train', *render*=False, *train_episodes*=1000, *test_episodes*=10, *max_steps*=200, *save_interval*=1000, *gamma*=0.99, *exploration_rate*=0.2, *exploration_final_eps*=0.01, *target_network_update_freq*=50, *batch_size*=32, *train_freq*=4, *learning_starts*=200, *plot_func*=None)

Parameters

- **env** – learning environment
- **mode** – train or test
- **render** – render each step
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **gamma** – reward decay factor
- **(float)** (*exploration_final_eps*) – fraction of entire training period over which the exploration rate is annealed
- **(float)** – final value of random action probability
- **(int)** (*learning_starts*) – update the target network every *target_network_update_freq* steps
- **(int)** – size of a batched sampled from replay buffer for training
- **(int)** – update the model every *train_freq* steps
- **(int)** – how many steps of the model to collect transitions for before learning starts
- **plot_func** – additional function for interactive module

load_ckpt (*env_name*)

load trained weights :return: None

save_ckpt (*env_name*)

save trained weights :return: None

store_transition (*s*, *a*, *r*, *s_*, *d*)

sync ()

Copy q network to target q network

update (*batch_size*, *gamma*)

5.3 Default Hyper-parameters

`rlzoo.algorithms.dqn.default.atari` (*env*, *default_seed*=False, ****kwargs**)

`rlzoo.algorithms.dqn.default.classic_control` (*env*, *default_seed*=False, ****kwargs**)

6.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import PG
4
5  AlgName = 'PG'
6  EnvName = 'PongNoFrameskip-v4'
7  EnvType = 'atari'
8
9  # EnvName = 'CartPole-v0'
10 # EnvType = 'classic_control'
11
12 # EnvName = 'BipedalWalker-v2'
13 # EnvType = 'box2d'
14
15 # EnvName = 'Ant-v2'
16 # EnvType = 'mujoco'
17
18 # EnvName = 'FetchPush-v1'
19 # EnvType = 'robotics'
20
21 # EnvName = 'FishSwim-v0'
22 # EnvType = 'dm_control'
23
24 # EnvName = 'ReachTarget'
25 # EnvType = 'rlbench'
26
27 env = build_env(EnvName, EnvType)
28 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
29 alg = eval(AlgName+'(**alg_params)')
30 alg.learn(env=env, mode='train', render=False, **learn_params)
31 alg.learn(env=env, mode='test', render=True, **learn_params)
```

6.2 Vanilla Policy Gradient

```
class rlzoo.algorithms.pg.pg.PG (net_list, optimizers_list)
    PG class

    get_action (s)
        choose action with probabilities.

        Parameters s – state

        Returns act

    get_action_greedy (s)
        choose action with greedy policy

        Parameters s – state

        Returns act

    learn (env, train_episodes=200, test_episodes=100, max_steps=200, save_interval=100,
        mode='train', render=False, gamma=0.95, plot_func=None)

        Parameters

        • env – learning environment

        • train_episodes – total number of episodes for training

        • test_episodes – total number of episodes for testing

        • max_steps – maximum number of steps for one episode

        • save_interval – time steps for saving

        • mode – train or test

        • render – render each step

        • gamma – reward decay

        • plot_func – additional function for interactive module

        Returns None

    load_ckpt (env_name)
        load trained weights

        Returns None

    save_ckpt (env_name)
        save trained weights

        Returns None

    store_transition (s, a, r)
        store data in memory buffer

        Parameters

        • s – state

        • a – act

        • r – reward

        Returns
```

update (*gamma*)
update policy parameters via stochastic gradient ascent
Returns None

6.3 Default Hyper-parameters

```
rlzoo.algorithms.pg.default.atari (env, default_seed=True)  
rlzoo.algorithms.pg.default.box2d (env, default_seed=True)  
rlzoo.algorithms.pg.default.classic_control (env, default_seed=True)  
rlzoo.algorithms.pg.default.dm_control (env, default_seed=True)  
rlzoo.algorithms.pg.default.mujoco (env, default_seed=True)  
rlzoo.algorithms.pg.default.rlbench (env, default_seed=True)  
rlzoo.algorithms.pg.default.robotics (env, default_seed=True)
```


7.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import AC
4
5  AlgName = 'AC'
6  EnvName = 'PongNoFrameskip-v4'
7  EnvType = 'atari'
8
9  # EnvName = 'Pendulum-v0'
10 # EnvType = 'classic_control'
11
12 # EnvName = 'BipedalWalker-v2'
13 # EnvType = 'box2d'
14
15 # EnvName = 'Ant-v2'
16 # EnvType = 'mujoco'
17
18 # EnvName = 'FetchPush-v1'
19 # EnvType = 'robotics'
20
21 # EnvName = 'FishSwim-v0'
22 # EnvType = 'dm_control'
23
24 # EnvName = 'ReachTarget'
25 # EnvType = 'rlbench'
26
27 env = build_env(EnvName, EnvType)
28 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
29 alg = eval(AlgName+'(**alg_params)')
30 alg.learn(env=env, mode='train', render=False, **learn_params)
31 alg.learn(env=env, mode='test', render=True, **learn_params)
```

7.2 Actor-Critic

```
class rlzoo.algorithms.ac.ac.AC (net_list, optimizers_list, gamma=0.9)
```

```
    get_action (s)
```

```
    get_action_greedy (s)
```

```
    learn (env, train_episodes=1000, test_episodes=500, max_steps=200, save_interval=100,  
         mode='train', render=False, plot_func=None)
```

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving the weights and plotting the results
- **mode** – ‘train’ or ‘test’
- **render** – if true, visualize the environment
- **plot_func** – additional function for interactive module

```
    load_ckpt (env_name)
```

```
    save_ckpt (env_name)
```

```
    update (s, a, r, s_)
```

7.3 Default Hyper-parameters

```
rlzoo.algorithms.ac.default.atari (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.box2d (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.classic_control (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.dm_control (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.mujoco (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.rlbench (env, default_seed=True)
```

```
rlzoo.algorithms.ac.default.robotics (env, default_seed=True)
```


8.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import A3C
4
5  AlgName = 'A3C'
6  EnvName = 'PongNoFrameskip-v4'
7  EnvType = 'atari'
8
9  # EnvName = 'Pendulum-v0' # only continuous action
10 # EnvType = 'classic_control'
11
12 # EnvName = 'BipedalWalker-v2'
13 # EnvType = 'box2d'
14
15 # EnvName = 'Ant-v2'
16 # EnvType = 'mujoco'
17
18 # EnvName = 'FetchPush-v1'
19 # EnvType = 'robotics'
20
21 # EnvName = 'FishSwim-v0'
22 # EnvType = 'dm_control'
23
24 number_workers = 2 # need to specify number of parallel workers
25 env = build_env(EnvName, EnvType, nenv=number_workers)
26 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
27 alg = eval(AlgName+'(**alg_params)')
28 alg.learn(env=env, mode='train', render=False, **learn_params)
29 alg.learn(env=env, mode='test', render=True, **learn_params)
```

8.2 Asynchronous Advantage Actor-Critic

```
class rlzoo.algorithms.a3c.a3c.A3C (net_list, optimizers_list, entropy_beta=0.005)
```

```
learn (env, train_episodes=1000, test_episodes=10, max_steps=150, render=False, n_workers=1, update_itr=10, gamma=0.99, save_interval=500, mode='train', plot_func=None)
```

Parameters

- **env** – a list of same learning environments
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **render** – render or not
- **n_workers** – manually set number of workers
- **update_itr** – update global policy after several episodes
- **gamma** – reward discount factor
- **save_interval** – timesteps for saving the weights and plotting the results
- **mode** – train or test
- **plot_func** – additional function for interactive module

8.3 Default Hyper-parameters

```
rlzoo.algorithms.a3c.default.atari (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.box2d (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.classic_control (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.dm_control (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.mujoco (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.rlbench (env, default_seed=True)
```

```
rlzoo.algorithms.a3c.default.robotics (env, default_seed=True)
```

9.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import DDPG
4
5  AlgName = 'DDPG'
6  EnvName = 'Pendulum-v0' # only continuous action
7  EnvType = 'classic_control'
8
9  # EnvName = 'BipedalWalker-v2'
10 # EnvType = 'box2d'
11
12 # EnvName = 'Ant-v2'
13 # EnvType = 'mujoco'
14
15 # EnvName = 'FetchPush-v1'
16 # EnvType = 'robotics'
17
18 # EnvName = 'FishSwim-v0'
19 # EnvType = 'dm_control'
20
21 # EnvName = 'ReachTarget'
22 # EnvType = 'rlbench'
23
24 env = build_env(EnvName, EnvType)
25 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
26 alg = eval(AlgName+'(**alg_params)')
27 alg.learn(env=env, mode='train', render=False, **learn_params)
28 alg.learn(env=env, mode='test', render=True, **learn_params)
```

9.2 Deep Deterministic Policy Gradient

```
class rlzoo.algorithms.ddpg.ddpg.DDPG(net_list, optimizers_list, replay_buffer_size, ac-  
                                     tion_range=1.0, tau=0.01)
```

DDPG class

```
ema_update()
```

Soft updating by exponential smoothing

Returns None

```
get_action(s, noise_scale)
```

Choose action with exploration

Parameters *s* – state

Returns action

```
get_action_greedy(s)
```

Choose action

Parameters *s* – state

Returns action

```
learn(env, train_episodes=200, test_episodes=100, max_steps=200, save_interval=10, ex-  
      plre_steps=500, mode='train', render=False, batch_size=32, gamma=0.9, noise_scale=1.0,  
      noise_scale_decay=0.995, plot_func=None)
```

learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **explore_steps** – for random action sampling in the beginning of training
- **mode** – train or test mode
- **render** – render each step
- **batch_size** – update batch size
- **gamma** – reward decay factor
- **noise_scale** – range of action noise for exploration
- **noise_scale_decay** – noise scale decay factor
- **plot_func** – additional function for interactive module

Returns None

```
load_ckpt(env_name)
```

load trained weights

Returns None

```
sample_action()
```

generate random actions for exploration

save_ckpt (*env_name*)
save trained weights

Returns None

store_transition (*s, a, r, s_, d*)
Store data in data buffer

Parameters

- **s** – state
- **a** – act
- **r** – reward
- **s** – next state

Returns None

update (*batch_size, gamma*)
Update parameters

Parameters

- **batch_size** – update batch size
- **gamma** – reward decay factor

Returns

9.3 Default Hyper-parameters

```
rlzoo.algorithms.ddpg.default.box2d (env, default_seed=True)
rlzoo.algorithms.ddpg.default.classic_control (env, default_seed=True)
rlzoo.algorithms.ddpg.default.dm_control (env, default_seed=True)
rlzoo.algorithms.ddpg.default.mujoco (env, default_seed=True)
rlzoo.algorithms.ddpg.default.rlbench (env, default_seed=True)
rlzoo.algorithms.ddpg.default.robotics (env, default_seed=True)
```


10.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import TD3
4
5  AlgName = 'TD3'
6  EnvName = 'Pendulum-v0' # only continuous action
7  EnvType = 'classic_control'
8
9  # EnvName = 'BipedalWalker-v2'
10 # EnvType = 'box2d'
11
12 # EnvName = 'Ant-v2'
13 # EnvType = 'mujoco'
14
15 # EnvName = 'FetchPush-v1'
16 # EnvType = 'robotics'
17
18 # EnvName = 'FishSwim-v0'
19 # EnvType = 'dm_control'
20
21 # EnvName = 'ReachTarget'
22 # EnvType = 'rlbench'
23
24 env = build_env(EnvName, EnvType)
25 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
26 alg = eval(AlgName+'(**alg_params)')
27 alg.learn(env=env, mode='train', render=False, **learn_params)
28 alg.learn(env=env, mode='test', render=True, **learn_params)
```

10.2 Twin Delayed DDPG

```
class rlzoo.algorithms.td3.td3.TD3 (net_list, optimizers_list, re-  
                                     play_buffer_capacity=500000.0, pol-  
                                     icy_target_update_interval=5)  
  
    twin-delayed ddpg  
  
    evaluate (state, eval_noise_scale, target=False)  
        generate action with state for calculating gradients;  
  
        Parameters eval_noise_scale – as the trick of target policy smoothing, for generating  
            noisy actions.  
  
    get_action (state, explore_noise_scale)  
        generate action with state for interaction with environment  
  
    get_action_greedy (state)  
        generate action with state for interaction with environment  
  
    learn (env, train_episodes=1000, test_episodes=1000, max_steps=150, batch_size=64, ex-  
        plore_steps=500, update_itr=3, reward_scale=1.0, save_interval=10, explore_noise_scale=1.0,  
        eval_noise_scale=0.5, mode='train', render=False, plot_func=None)  
  
        Parameters  
  
        • env – learning environment  
  
        • train_episodes – total number of episodes for training  
  
        • test_episodes – total number of episodes for testing  
  
        • max_steps – maximum number of steps for one episode  
  
        • batch_size – update batchsize  
  
        • explore_steps – for random action sampling in the beginning of training  
  
        • update_itr – repeated updates for single step  
  
        • reward_scale – value range of reward  
  
        • save_interval – timesteps for saving the weights and plotting the results  
  
        • explore_noise_scale – range of action noise for exploration  
  
        • eval_noise_scale – range of action noise for evaluation of action value  
  
        • mode – ‘train’ or ‘test’  
  
        • render – if true, visualize the environment  
  
        • plot_func – additional function for interactive module  
  
    load_ckpt (env_name)  
  
    sample_action ()  
        generate random actions for exploration  
  
    save_ckpt (env_name)  
  
    target_ini (net, target_net)  
        hard-copy update for initializing target networks  
  
    target_soft_update (net, target_net, soft_tau)  
        soft update the target net with Polyak averaging
```


update (*batch_size*, *eval_noise_scale*, *reward_scale*=1.0, *gamma*=0.9, *soft_tau*=0.01)
update all networks in TD3

10.3 Default Hyper-parameters

```
rlzoo.algorithms.td3.default.box2d (env, default_seed=True)  
rlzoo.algorithms.td3.default.classic_control (env, default_seed=True)  
rlzoo.algorithms.td3.default.dm_control (env, default_seed=True)  
rlzoo.algorithms.td3.default.mujoco (env, default_seed=True)  
rlzoo.algorithms.td3.default.rlbench (env, default_seed=True)  
rlzoo.algorithms.td3.default.robotics (env, default_seed=True)
```


11.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import SAC
4
5  AlgName = 'SAC'
6  EnvName = 'Pendulum-v0' # only continuous action
7  EnvType = 'classic_control'
8
9  # EnvName = 'BipedalWalker-v2'
10 # EnvType = 'box2d'
11
12 # EnvName = 'Ant-v2'
13 # EnvType = 'mujoco'
14
15 # EnvName = 'FetchPush-v1'
16 # EnvType = 'robotics'
17
18 # EnvName = 'FishSwim-v0'
19 # EnvType = 'dm_control'
20
21 # EnvName = 'ReachTarget'
22 # EnvType = 'rlbench'
23
24 env = build_env(EnvName, EnvType)
25 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
26 alg = eval(AlgName+'(**alg_params)')
27 alg.learn(env=env, mode='train', render=False, **learn_params)
28 alg.learn(env=env, mode='test', render=True, **learn_params)
```

11.2 Soft Actor-Critic

```
class rlzoo.algorithms.sac.sac.SAC (net_list, optimizers_list, re-  
                                     play_buffer_capacity=500000.0)  
    Soft Actor-Critic  
  
    evaluate (state, epsilon=1e-06)  
        generate action with state for calculating gradients  
  
    get_action (state)  
        generate action with state for interaction with environment  
  
    get_action_greedy (state)  
        generate action with state for interaction with environment  
  
    learn (env, train_episodes=1000, test_episodes=1000, max_steps=150, batch_size=64, ex-  
        ploration_steps=500, update_itr=3, policy_target_update_interval=3, reward_scale=1.0,  
        save_interval=20, mode='train', AUTO_ENTROPY=True, render=False, plot_func=None)  
        Parameters  
        • env – learning environment  
        • train_episodes – total number of episodes for training  
        • test_episodes – total number of episodes for testing  
        • max_steps – maximum number of steps for one episode  
        • batch_size – update batchsize  
        • explore_steps – for random action sampling in the beginning of training  
        • update_itr – repeated updates for single step  
        • policy_target_update_interval – delayed update for the policy network and  
          target networks  
        • reward_scale – value range of reward  
        • save_interval – timesteps for saving the weights and plotting the results  
        • mode – ‘train’ or ‘test’  
        • AUTO_ENTROPY – automatically updating variable alpha for entropy  
        • render – if true, visualize the environment  
        • plot_func – additional function for interactive module  
  
    load_ckpt (env_name)  
        load trained weights  
  
    sample_action ()  
        generate random actions for exploration  
  
    save_ckpt (env_name)  
        save trained weights  
  
    target_ini (net, target_net)  
        hard-copy update for initializing target networks  
  
    target_soft_update (net, target_net, soft_tau)  
        soft update the target net with Polyak averaging
```

```
update (batch_size, reward_scale=10.0, auto_entropy=True, target_entropy=-2, gamma=0.99,  
        soft_tau=0.01)  
update all networks in SAC
```

11.3 Default Hyper-parameters

```
rlzoo.algorithms.sac.default.box2d (env, default_seed=True)  
rlzoo.algorithms.sac.default.classic_control (env, default_seed=True)  
rlzoo.algorithms.sac.default.dm_control (env, default_seed=True)  
rlzoo.algorithms.sac.default.mujoco (env, default_seed=True)  
rlzoo.algorithms.sac.default.rlbench (env, default_seed=True)  
rlzoo.algorithms.sac.default.robotics (env, default_seed=True)
```


12.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import TD3
4
5  AlgName = 'TRPO'
6  EnvName = 'PongNoFrameskip-v4'
7  EnvType = 'atari'
8
9  # EnvName = 'CartPole-v0'
10 # EnvType = 'classic_control'
11
12 # EnvName = 'BipedalWalker-v2'
13 # EnvType = 'box2d'
14
15 # EnvName = 'Ant-v2'
16 # EnvType = 'mujoco'
17
18 # EnvName = 'FetchPush-v1'
19 # EnvType = 'robotics'
20
21 # EnvName = 'FishSwim-v0'
22 # EnvType = 'dm_control'
23
24 # EnvName = 'ReachTarget'
25 # EnvType = 'rlbench'
26
27 env = build_env(EnvName, EnvType)
28 alg_params, learn_params = call_default_params(env, EnvType, AlgName)
29 alg = eval(AlgName+'(**alg_params)')
30 alg.learn(env=env, mode='train', render=False, **learn_params)
31 alg.learn(env=env, mode='test', render=True, **learn_params)
```

12.2 Trust Region Policy Optimization

```
class rlzoo.algorithms.trpo.trpo.TRPO (net_list, optimizers_list, damping_coeff=0.1,  
                                         cg_iters=10, delta=0.01)  
    trpo class  
  
    a_train (s, a, adv, oldpi_prob, backtrack_iters, backtrack_coeff)  
  
    static assign_params_from_flat (x, params)  
        assign params from flat input  
  
        Parameters  
        • x –  
        • params –  
  
        Returns group  
  
    c_train (tfdc_r, s)  
        Update actor network  
  
        Parameters  
        • tfdc_r – cumulative reward  
        • s – state  
  
        Returns None  
  
    cal_adv (tfs, tfdc_r)  
        Calculate advantage  
  
        Parameters  
        • tfs – state  
        • tfdc_r – cumulative reward  
  
        Returns advantage  
  
    cg (Ax, b)  
        Conjugate gradient algorithm (see https://en.wikipedia.org/wiki/Conjugate\_gradient\_method)  
  
    eval (bs, ba, badv, oldpi_prob)  
  
    static flat_concat (xs)  
        flat concat input  
  
        Parameters xs – a list of tensor  
  
        Returns flat tensor  
  
    get_action (s)  
        Choose action  
  
        Parameters s – state  
  
        Returns clipped act  
  
    get_action_greedy (s)  
        Choose action  
  
        Parameters s – state  
  
        Returns clipped act
```


get_pi_params ()
get actor trainable parameters

Returns flat actor trainable parameters

get_v (s)
Compute value

Parameters s – state

Returns value

gradient (inputs)
pi gradients

Parameters inputs – a list of x_ph, a_ph, adv_ph, ret_ph, logp_old_ph and other inputs

Returns gradient

hessian_vector_product (s, a, adv, oldpi_prob, v_ph)

learn (env, train_episodes=200, test_episodes=100, max_steps=200, save_interval=10, gamma=0.9, mode='train', render=False, batch_size=32, backtrack_iters=10, backtrack_coeff=0.8, train_critic_iters=80, plot_func=None)
learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **gamma** – reward discount factor
- **mode** – train or test
- **render** – render each step
- **batch_size** – update batch size
- **backtrack_iters** – Maximum number of steps allowed in the backtracking line search
- **backtrack_coeff** – How far back to step during backtracking line search
- **train_critic_iters** – critic update iteration steps

Returns None

load_ckpt (env_name)
load trained weights

Returns None

pi_loss (inputs)
calculate pi loss

Parameters inputs – a list of x_ph, a_ph, adv_ph, ret_ph, logp_old_ph and other inputs

Returns pi loss

save_ckpt (env_name)
save trained weights

Returns None

set_pi_params (*v_ph*)
set actor trainable parameters

Parameters *v_ph* – inputs

Returns None

update (*bs, ba, br, train_critic_iters, backtrack_iters, backtrack_coeff*)
update trpo

Returns None

12.3 Default Hyper-parameters

```
rlzoo.algorithms.trpo.default.atari (env, default_seed=True)  
rlzoo.algorithms.trpo.default.box2d (env, default_seed=True)  
rlzoo.algorithms.trpo.default.classic_control (env, default_seed=True)  
rlzoo.algorithms.trpo.default.dm_control (env, default_seed=True)  
rlzoo.algorithms.trpo.default.mujoco (env, default_seed=True)  
rlzoo.algorithms.trpo.default.rlbench (env, default_seed=True)  
rlzoo.algorithms.trpo.default.robotics (env, default_seed=True)
```

13.1 Example

```
1  from rlzoo.common.env_wrappers import build_env
2  from rlzoo.common.utils import call_default_params
3  from rlzoo.algorithms import PPO
4
5  EnvName = 'PongNoFrameskip-v4'
6  EnvType = 'atari'
7
8  # EnvName = 'Pendulum-v0'
9  # EnvType = 'classic_control'
10
11 # EnvName = 'BipedalWalker-v2'
12 # EnvType = 'box2d'
13
14 # EnvName = 'Ant-v2'
15 # EnvType = 'mujoco'
16
17 # EnvName = 'FetchPush-v1'
18 # EnvType = 'robotics'
19
20 # EnvName = 'FishSwim-v0'
21 # EnvType = 'dm_control'
22
23 # EnvName = 'ReachTarget'
24 # EnvType = 'rlbench'
25
26 env = build_env(EnvName, EnvType)
27 alg_params, learn_params = call_default_params(env, EnvType, 'PPO')
28 alg = PPO(method='clip', **alg_params) # specify 'clip' or 'penalty' method for PPO
29 alg.learn(env=env, mode='train', render=False, **learn_params)
30 alg.learn(env=env, mode='test', render=False, **learn_params)
```

13.2 Proximal Policy Optimization (Penalty)

```
class rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY (net_list,           op-  
                                                         timizers_list,  
                                                         kl_target=0.01,  
                                                         lam=0.5)
```

PPO class

```
a_train (tfs, tfa, tfadv, oldpi_prob)  
    Update policy network
```

Parameters

- **tfs** – state
- **tfa** – act
- **tfadv** – advantage

Returns

```
c_train (tfdc_r, s)  
    Update actor network
```

Parameters

- **tfdc_r** – cumulative reward
- **s** – state

Returns None

```
cal_adv (tfs, tfdc_r)  
    Calculate advantage
```

Parameters

- **tfs** – state
- **tfdc_r** – cumulative reward

Returns advantage

```
get_action (s)  
    Choose action
```

Parameters **s** – state

Returns clipped act

```
get_action_greedy (s)  
    Choose action
```

Parameters **s** – state

Returns clipped act

```
get_v (s)  
    Compute value
```

Parameters **s** – state

Returns value

learn (*env*, *train_episodes=1000*, *test_episodes=10*, *max_steps=200*, *save_interval=10*, *gamma=0.9*,
mode='train', *render=False*, *batch_size=32*, *a_update_steps=10*, *c_update_steps=10*,
plot_func=None)
 learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **gamma** – reward discount factor
- **mode** – train or test
- **render** – render each step
- **batch_size** – update batch size
- **a_update_steps** – actor update iteration steps
- **c_update_steps** – critic update iteration steps
- **plot_func** – additional function for interactive module

Returns None

load_ckpt (*env_name*)
 load trained weights

Returns None

save_ckpt (*env_name*)
 save trained weights

Returns None

update (*s*, *a*, *r*, *a_update_steps*, *c_update_steps*)
 Update parameter with the constraint of KL divergent

Parameters

- **s** – state
- **a** – act
- **r** – reward

Returns None

13.3 Proximal Policy Optimization (Clip)

class rlzoo.algorithms.ppo_clip.ppo_clip.**PPO_CLIP** (*net_list*, *optimizers_list*, *epsilon=0.2*)

PPO class

a_train (*tfs*, *tfa*, *tfadv*, *oldpi_prob*)
 Update policy network

Parameters

- **tfs** – state
- **tfa** – act
- **tfadv** – advantage
- **oldpi_prob** – old policy distribution

Returns None

c_train (*tfdc_r*, *s*)
Update actor network

Parameters

- **tfdc_r** – cumulative reward
- **s** – state

Returns None

cal_adv (*tfs*, *tfdc_r*)
Calculate advantage

Parameters

- **tfs** – state
- **tfdc_r** – cumulative reward

Returns advantage

get_action (*s*)
Choose action

Parameters **s** – state

Returns clipped act

get_action_greedy (*s*)
Choose action :param s: state :return: clipped act

get_v (*s*)
Compute value

Parameters **s** – state

Returns value

learn (*env*, *train_episodes=200*, *test_episodes=100*, *max_steps=200*, *save_interval=10*, *gamma=0.9*,
mode='train', *render=False*, *batch_size=32*, *a_update_steps=10*, *c_update_steps=10*,
plot_func=None)
learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – timesteps for saving
- **gamma** – reward discount factor
- **mode** – train or test

- **render** – render each step
- **batch_size** – update batchsize
- **a_update_steps** – actor update iteration steps
- **c_update_steps** – critic update iteration steps
- **plot_func** – additional function for interactive module

Returns None

load_ckpt (*env_name*)
load trained weights

Returns None

save_ckpt (*env_name*)
save trained weights

Returns None

update (*s, a, r, a_update_steps, c_update_steps*)
Update parameter with the constraint of KL divergent

Parameters

- **s** – state
- **a** – act
- **r** – reward

Returns None

13.4 Default Hyper-parameters

```
rlzoo.algorithms.ppo.default.atari (env, default_seed=True)
rlzoo.algorithms.ppo.default.box2d (env, default_seed=True)
rlzoo.algorithms.ppo.default.classic_control (env, default_seed=True)
rlzoo.algorithms.ppo.default.dm_control (env, default_seed=True)
rlzoo.algorithms.ppo.default.mujoco (env, default_seed=True)
rlzoo.algorithms.ppo.default.rlbench (env, default_seed=True)
rlzoo.algorithms.ppo.default.robotics (env, default_seed=True)
```


14.1 Example

```
1 from rlzoo.common.env_wrappers import build_env
2 from rlzoo.common.utils import call_default_params
3 from rlzoo.algorithms import DPPO
4
5 EnvName = 'PongNoFrameskip-v4'
6 EnvType = 'atari'
7
8 # EnvName = 'Pendulum-v0'
9 # EnvType = 'classic_control'
10
11 # EnvName = 'BipedalWalker-v2'
12 # EnvType = 'box2d'
13
14 # EnvName = 'Ant-v2'
15 # EnvType = 'mujoco'
16
17 # EnvName = 'FetchPush-v1'
18 # EnvType = 'robotics'
19
20 # EnvName = 'FishSwim-v0'
21 # EnvType = 'dm_control'
22
23 # EnvName = 'ReachTarget'
24 # EnvType = 'rlbench'
25
26 number_workers = 2 # need to specify number of parallel workers
27 env = build_env(EnvName, EnvType, nenv=number_workers)
28 alg_params, learn_params = call_default_params(env, EnvType, 'DPPO')
29 alg = DPPO(method='penalty', **alg_params) # specify 'clip' or 'penalty' method for
↪ PPO
30 alg.learn(env=env, mode='train', render=False, **learn_params)
```

(continues on next page)

(continued from previous page)

```
31 alg.learn(env=env, mode='test', render=True, **learn_params)
```

14.2 Distributed Proximal Policy Optimization (Penalty)

```
class rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY(net_list, op-  
timizers_list,  
kl_target=0.01,  
lam=0.5)
```

PPO class

a_train (*tfs, tfa, tfadv, oldpi_prob*)
Update policy network

Parameters

- **tfs** – state
- **tfa** – act
- **tfadv** – advantage

Returns

c_train (*tfdc_r, s*)
Update actor network

Parameters

- **tfdc_r** – cumulative reward
- **s** – state

Returns None

cal_adv (*tfs, tfdc_r*)
Calculate advantage

Parameters

- **tfs** – state
- **tfdc_r** – cumulative reward

Returns advantage

get_action (*s*)
Choose action

Parameters **s** – state

Returns clipped act

get_action_greedy (*s*)
Choose action

Parameters **s** – state

Returns clipped act

get_v (*s*)
Compute value

Parameters **s** – state

Returns value

learn (*env*, *train_episodes=200*, *test_episodes=100*, *max_steps=200*, *save_interval=10*, *gamma=0.9*,
mode='train', *render=False*, *batch_size=32*, *a_update_steps=10*, *c_update_steps=10*,
n_workers=4, *plot_func=None*)
 learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **gamma** – reward discount factor
- **mode** – train or test
- **render** – render each step
- **batch_size** – update batch size
- **a_update_steps** – actor update iteration steps
- **c_update_steps** – critic update iteration steps
- **n_workers** – number of workers
- **plot_func** – additional function for interactive module

Returns None

load_ckpt (*env_name*)
 load trained weights

Returns None

save_ckpt (*env_name*)
 save trained weights

Returns None

update (*a_update_steps*, *c_update_steps*, *save_interval*, *env*)
 Update

Parameters

- **a_update_steps** – actor update steps
- **c_update_steps** – critic update steps
- **save_interval** – save interval
- **env** – environment

Returns None

14.3 Distributed Proximal Policy Optimization (Clip)

```
class rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP(net_list, optimizers_list, epsilon=0.2)
```

PPO class

```
a_train(tfs, tfa, tfadv, oldpi_prob)
```

Update policy network

Parameters

- **tfs** – state
- **tfa** – act
- **tfadv** – advantage
- **oldpi_prob** – old pi probability of a in s

Returns

```
c_train(tfdc_r, s)
```

Update actor network

Parameters

- **tfdc_r** – cumulative reward
- **s** – state

Returns None

```
cal_adv(tfs, tfdc_r)
```

Calculate advantage

Parameters

- **tfs** – state
- **tfdc_r** – cumulative reward

Returns advantage

```
get_action(s)
```

Choose action

Parameters **s** – state

Returns clipped act

```
get_action_greedy(s)
```

Choose action

Parameters **s** – state

Returns clipped act

```
get_v(s)
```

Compute value

Parameters **s** – state

Returns value

learn (*env*, *train_episodes=200*, *test_episodes=100*, *max_steps=200*, *save_interval=10*, *gamma=0.9*,
mode='train', *render=False*, *batch_size=32*, *a_update_steps=10*, *c_update_steps=10*,
n_workers=4, *plot_func=None*)
 learn function

Parameters

- **env** – learning environment
- **train_episodes** – total number of episodes for training
- **test_episodes** – total number of episodes for testing
- **max_steps** – maximum number of steps for one episode
- **save_interval** – time steps for saving
- **gamma** – reward discount factor
- **mode** – train or test
- **render** – render each step
- **batch_size** – update batch size
- **a_update_steps** – actor update iteration steps
- **c_update_steps** – critic update iteration steps
- **n_workers** – number of workers
- **plot_func** – additional function for interactive module

Returns None

load_ckpt (*env_name*)
 load trained weights

Returns None

save_ckpt (*env_name*)
 save trained weights

Returns None

update (*a_update_steps*, *c_update_steps*, *save_interval*, *env*)
 Update

Parameters

- **a_update_steps** – actor update steps
- **c_update_steps** – critic update steps
- **save_interval** – save interval
- **env** – environment

Returns None

14.4 Default Hyper-parameters

`rlzoo.algorithms.dppo.default.atari` (*env*, *default_seed=True*)

`rlzoo.algorithms.dppo.default.box2d` (*env*, *default_seed=True*)

```
rlzoo.algorithms.dppo.default.classic_control (env, default_seed=True)
rlzoo.algorithms.dppo.default.dm_control (env, default_seed=True)
rlzoo.algorithms.dppo.default.mujoco (env, default_seed=True)
rlzoo.algorithms.dppo.default.robotics (env, default_seed=True)
```

15.1 Basic Networks in RLzoo

Basic neural networks

`rlzoo.common.basic_nets.CNN` (*input_shape*, *conv_kwargs=None*)

Multiple convolutional layers for approximation Default setting is equal to architecture used in DQN

Parameters

- **input_shape** – (tuple[int]) (H, W, C)
- **conv_kwargs** – (list[param]) list of conv parameters for `tl.layers.Conv2d`

Return: input tensor, output tensor

`rlzoo.common.basic_nets.CNNModel` (*input_shape*, *conv_kwargs=None*)

Multiple convolutional layers for approximation Default setting is equal to architecture used in DQN

Parameters

- **input_shape** – (tuple[int]) (H, W, C)
- **conv_kwargs** – (list[param]) list of conv parameters for `tl.layers.Conv2d`

Return: `tl.model.Model`

`rlzoo.common.basic_nets.CreateInputLayer` (*state_space*, *conv_kwargs=None*)

`rlzoo.common.basic_nets.MLP` (*input_dim*, *hidden_dim_list*, *w_init=<tensorflow.python.ops.init_ops_v2.Orthogonal object>*, *activation=<function relu>*, **args*, ***kwargs*)

Multiple fully-connected layers for approximation

Parameters

- **input_dim** – (int) size of input tensor
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers

- **w_init** – (callable) initialization method for weights
- **activation** – (callable) activation function of hidden layers

Return: input tensor, output tensor

```
rlzoo.common.basic_nets.MLPModel(input_dim, hidden_dim_list,  
                                  w_init=<tensorflow.python.ops.init_ops_v2.Orthogonal  
                                  object>, activation=<function relu>, *args, **kwargs)
```

Multiple fully-connected layers for approximation

Parameters

- **input_dim** – (int) size of input tensor
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) initialization method for weights
- **activation** – (callable) activation function of hidden layers

Return: input tensor, output tensor

16.1 Policy Networks in RLzoo

```
class rlzoo.common.policy_networks.StochasticContinuousPolicyNetwork (state_shape,
                                                                    ac-
                                                                    tion_shape,
                                                                    hid-
                                                                    den_dim_list,
                                                                    w_init=<tensorflow.python.ops.
                                                                    ob-
                                                                    ject>,
                                                                    activa-
                                                                    tion=<function
                                                                    relu>,
                                                                    out-
                                                                    put_activation=None,
                                                                    log_std_min=-
                                                                    20,
                                                                    log_std_max=2,
                                                                    train-
                                                                    able=True)
```

Bases: tensorflow.models.core.Model

```
__init__ (state_shape, action_shape, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
object>, activation=<function relu>, output_activation=None, log_std_min=-20,
log_std_max=2, trainable=True)
```

Stochastic continuous policy network with multiple fully-connected layers or convolutional layers (accord-
ing to state shape)

Parameters

- **state_shape** – (tuple[int]) shape of the state, for example, (state_dim,) for single-dimensional state

- **action_shape** – (tuple[int]) shape of the action, for example, (action_dim,) for single-dimensional action
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **log_std_min** – (float) lower bound of standard deviation of action
- **log_std_max** – (float) upper bound of standard deviation of action
- **trainable** – (bool) set training and evaluation mode

```
__module__ = 'rlzoo.common.policy_networks'
```

```
class rlzoo.common.policy_networks.DeterministicContinuousPolicyNetwork (state_shape,
                                                                           ac-
                                                                           tion_shape,
                                                                           hid-
                                                                           den_dim_list,
                                                                           w_init=<tensorflow.python.
                                                                           ob-
                                                                           ject>,
                                                                           ac-
                                                                           ti-
                                                                           va-
                                                                           tion=<function
                                                                           relu>,
                                                                           out-
                                                                           put_activation=<function
                                                                           tanh>,
                                                                           train-
                                                                           able=True)
```

Bases: `tensorlayer.models.core.Model`

```
__init__ (state_shape, action_shape, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
           object>, activation=<function relu>, output_activation=<function tanh>, trainable=True)
Deterministic continuous policy network with multiple fully-connected layers or convolutional layers (ac-
cording to state shape)
```

Parameters

- **state_shape** – (tuple[int]) shape of the state, for example, (state_dim,) for single-dimensional state
- **action_shape** – (tuple[int]) shape of the action, for example, (action_dim,) for single-dimensional action
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **trainable** – (bool) set training and evaluation mode

```
__module__ = 'rlzoo.common.policy_networks'
```

```

class rlzoo.common.policy_networks.DeterministicPolicyNetwork (state_space,
                                                                action_space,
                                                                hidden_dim_list,
                                                                w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
                                                                object>, activation=<function
                                                                relu>, output_activation=<function
                                                                tanh>, trainable=True,
                                                                name=None)

```

Bases: tensorflow.models.core.Model

__call__ (states, *args, **kwargs)

Forward input tensors through this network by calling.

inputs [Tensor or list of Tensors, numpy.ndarray of list of numpy.ndarray] Inputs for network forwarding

is_train [boolean] Network's mode for this time forwarding. If 'is_train' == True, this network is set as training mode. If 'is_train' == False, this network is set as evaluation mode

kwargs : For other keyword-only arguments.

__init__ (state_space, action_space, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal object>, activation=<function relu>, output_activation=<function tanh>, trainable=True, name=None)

Deterministic continuous/discrete policy network with multiple fully-connected layers

Parameters

- **state_space** – (gym.spaces) space of the state from gym environments
- **action_space** – (gym.spaces) space of the action from gym environments
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **trainable** – (bool) set training and evaluation mode

__module__ = 'rlzoo.common.policy_networks'

action_shape

action_space

random_sample ()

generate random actions for exploration

state_shape

state_space

```
class rlzoo.common.policy_networks.StochasticPolicyNetwork (state_space,      ac-
                                                           tion_space,      hid-
                                                           den_dim_list,
                                                           w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
                                                           object>,      activa-
                                                           tion=<function
                                                           relu>,      out-
                                                           put_activation=<function
                                                           tanh>, log_std_min=-
                                                           20, log_std_max=2,
                                                           trainable=True,
                                                           name=None,
                                                           state_conditioned=False)
```

Bases: tensorflow.models.core.Model

__call__ (states, *args, greedy=False, **kwargs)

Forward input tensors through this network by calling.

inputs [Tensor or list of Tensors, numpy.ndarray of list of numpy.ndarray] Inputs for network forwarding

is_train [boolean] Network's mode for this time forwarding. If 'is_train' == True, this network is set as training mode. If 'is_train' == False, this network is set as evaluation mode

kwargs : For other keyword-only arguments.

__init__ (state_space, action_space, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal object>, activation=<function relu>, output_activation=<function tanh>, log_std_min=-20, log_std_max=2, trainable=True, name=None, state_conditioned=False)

Stochastic continuous/discrete policy network with multiple fully-connected layers

Parameters

- **state_space** – (gym.spaces) space of the state from gym environments
- **action_space** – (gym.spaces) space of the action from gym environments
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **log_std_min** – (float) lower bound of standard deviation of action
- **log_std_max** – (float) upper bound of standard deviation of action
- **trainable** – (bool) set training and evaluation mode

Tips: We recommend to use `tf.nn.tanh` for `output_activation`, especially for continuous action space, to ensure the final action range is exactly the same as declared in action space after action normalization.

__module__ = 'rlzoo.common.policy_networks'

action_shape

action_space

random_sample ()

generate random actions for exploration

state_shape

state_space

17.1 Value Networks in RLzoo

```
class rlzoo.common.value_networks.ValueNetwork (state_space,             hidden_dim_list,
                                                w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
                                                object>, activation=<function relu>,
                                                output_activation=None,         trainable=True, name=None)
```

```
__call__ (states, *args, **kwargs)
```

Forward input tensors through this network by calling.

inputs [Tensor or list of Tensors, numpy.ndarray of list of numpy.ndarray] Inputs for network forwarding

is_train [boolean] Network's mode for this time forwarding. If 'is_train' == True, this network is set as training mode. If 'is_train' == False, this network is set as evaluation mode

kwargs: For other keyword-only arguments.

```
__init__ (state_space, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
        object>, activation=<function relu>, output_activation=None, trainable=True,
        name=None)
```

Value network with multiple fully-connected layers or convolutional layers (according to state shape)

Parameters

- **state_space** – (gym.spaces) space of the state from gym environments
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **trainable** – (bool) set training and evaluation mode

```
__module__ = 'rlzoo.common.value_networks'
```

state_shape

state_space

```
class rlzoo.common.value_networks.MlpQNetwork (state_shape,          ac-
                                              tion_shape,          hidden_dim_list,
                                              w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
                                              object>, activation=<function relu>, out-
                                              put_activation=None, trainable=True)
```

```
__init__ (state_shape, action_shape, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
object>, activation=<function relu>, output_activation=None, trainable=True)
```

Q-value network with multiple fully-connected layers

Inputs: (state tensor, action tensor)

Parameters

- **state_shape** – (tuple[int]) shape of the state, for example, (state_dim,) for single-dimensional state
- **action_shape** – (tuple[int]) shape of the action, for example, (action_dim,) for single-dimensional action
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers
- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **trainable** – (bool) set training and evaluation mode

```
__module__ = 'rlzoo.common.value_networks'
```

```
class rlzoo.common.value_networks.QNetwork (state_space, action_space, hidden_dim_list,
                                              w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
                                              object>, activation=<function relu>, out-
                                              put_activation=None, trainable=True,
                                              name=None, state_only=False, duel-
                                              ing=False)
```

```
__call__ (inputs, *args, **kwargs)
```

Forward input tensors through this network by calling.

inputs [Tensor or list of Tensors, numpy.ndarray of list of numpy.ndarray] Inputs for network forwarding

is_train [boolean] Network's mode for this time forwarding. If 'is_train' == True, this network is set as training mode. If 'is_train' == False, this network is set as evaluation mode

kwargs : For other keyword-only arguments.

```
__init__ (state_space, action_space, hidden_dim_list, w_init=<tensorflow.python.ops.init_ops_v2.GlorotNormal
object>, activation=<function relu>, output_activation=None, trainable=True,
name=None, state_only=False, dueling=False)
```

Q-value network with multiple fully-connected layers or convolutional layers (according to state shape)

Parameters

- **state_space** – (gym.spaces) space of the state from gym environments
- **action_space** – (gym.spaces) space of the action from gym environments
- **hidden_dim_list** – (list[int]) a list of dimensions of hidden layers

- **w_init** – (callable) weights initialization
- **activation** – (callable) activation function
- **output_activation** – (callable or None) output activation function
- **trainable** – (bool) set training and evaluation mode
- **name** – (str) name the model
- **state_only** – (bool) only input state or not, available in discrete action space
- **dueling** – (bool) whether use the dueling output or not, available in discrete action space

```
__module__ = 'rlzoo.common.value_networks'
```

```
action_shape
```

```
action_space
```

```
state_shape
```

```
state_space
```


18.1 Replay Buffer in RLzoo

Functions for utilization.

Requirements tensorflow==2.0.0a0 tensorlayer==2.0.1

```
class rlzoo.common.buffer.HindsightReplayBuffer(capacity, hindsight_freq, goal_type, re-  
                                              ward_func, done_func)
```

Bases: *rlzoo.common.buffer.ReplayBuffer*

Hindsight Experience Replay In this buffer, state is a tuple consists of (observation, goal)

GOAL_EPISODE = 'episode'

GOAL_FUTURE = 'future'

GOAL_RANDOM = 'random'

```
__init__(capacity, hindsight_freq, goal_type, reward_func, done_func)
```

Parameters

- **(int)** (*hindsight_freq*) – How many hindsight transitions will be generated for each real transition
- **(str)** (*goal_type*) – The generation method of hindsight goals. Should be HER_GOAL_*
- **(callable)** (*done_func*) – goal (np.array) X next_state (np.array) -> reward (float)
- **(callable)** – goal (np.array) X next_state (np.array) -> done_flag (bool)

```
__module__ = 'rlzoo.common.buffer'
```

```
push(*args, **kwargs)
```

```
push_episode(states, actions, rewards, next_states, dones)
```

```
class rlzoo.common.buffer.MinSegmentTree(capacity)
```

Bases: `rlzoo.common.buffer.SegmentTree`

```
__init__(capacity)
```

Build a Segment Tree data structure.

https://en.wikipedia.org/wiki/Segment_tree

Can be used as regular array, but with two important differences:

- setting item's value is slightly slower. It is $O(\lg \text{capacity})$ instead of $O(1)$.
- user has access to an efficient ($O(\log \text{segment size})$) *reduce* operation which reduces *operation* over a contiguous subsequence of items in the array.

Parameters

- capacity** – (int) Total size of the array - must be a power of two.
- operation** – (lambda obj, obj -> obj) and operation for combining elements (eg. sum, max) must form a mathematical group together with the set of possible values for array elements (i.e. be associative)
- neutral_element** – (obj) neutral element for the operation above. eg. float('-inf') for max and 0 for sum.

```
__module__ = 'rlzoo.common.buffer'
```

```
min(start=0, end=None)
```

Returns min(arr[start], ..., arr[end])

```
class rlzoo.common.buffer.PrioritizedReplayBuffer(capacity, alpha, beta)
```

Bases: `rlzoo.common.buffer.ReplayBuffer`

```
__init__(capacity, alpha, beta)
```

Create Prioritized Replay buffer.

Parameters

- capacity** – (int) Max number of transitions to store in the buffer. When the buffer overflows the old memories are dropped.
- alpha** – (float) how much prioritization is used (0 - no prioritization, 1 - full prioritization)

See Also: `ReplayBuffer.__init__`

```
__module__ = 'rlzoo.common.buffer'
```

```
push(*args)
```

See `ReplayBuffer.store_effect`

```
sample(batch_size)
```

Sample a batch of experiences

```
update_priorities(idxs, priorities)
```

Update priorities of sampled transitions

```
class rlzoo.common.buffer.ReplayBuffer(capacity)
```

Bases: `object`

A standard ring buffer for storing transitions and sampling for training

```
__dict__ = mappingproxy({'__module__': 'rlzoo.common.buffer', '__doc__': 'A standard
```

```

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

__len__()

__module__ = 'rlzoo.common.buffer'

__weakref__
    list of weak references to the object (if defined)

push(state, action, reward, next_state, done)

sample(batch_size)

class rlzoo.common.buffer.SegmentTree(capacity, operation, neutral_element)
    Bases: object

    __dict__ = mappingproxy({'__module__': 'rlzoo.common.buffer', '__init__': <function ...
    __getitem__(idx)

    __init__(capacity, operation, neutral_element)
        Build a Segment Tree data structure.

        https://en.wikipedia.org/wiki/Segment\_tree

        Can be used as regular array, but with two important differences:

        a) setting item's value is slightly slower. It is O(lg capacity) instead of O(1).

        b) user has access to an efficient ( O(log segment size) ) reduce operation which reduces operation over
           a contiguous subsequence of items in the array.

    Parameters

        • apacity – (int) Total size of the array - must be a power of two.

        • operation – (lambda obj, obj -> obj) and operation for combining elements (eg. sum,
          max) must form a mathematical group together with the set of possible values for array
          elements (i.e. be associative)

        • neutral_element – (obj) neutral element for the operation above. eg. float('-inf') for
          max and 0 for sum.

    __module__ = 'rlzoo.common.buffer'

    __setitem__(idx, val)

    __weakref__
        list of weak references to the object (if defined)

    reduce(start=0, end=None)
        Returns result of applying self.operation to a contiguous subsequence of the array.

    Parameters

        • start – (int) beginning of the subsequence

        • end – (int) end of the subsequences

    Returns: reduced: (obj) result of reducing self.operation over the specified range of array.

class rlzoo.common.buffer.SumSegmentTree(capacity)
    Bases: rlzoo.common.buffer.SegmentTree

```

`__init__(capacity)`

Build a Segment Tree data structure.

https://en.wikipedia.org/wiki/Segment_tree

Can be used as regular array, but with two important differences:

- setting item's value is slightly slower. It is $O(\lg \text{capacity})$ instead of $O(1)$.
- user has access to an efficient ($O(\log \text{segment size})$) *reduce* operation which reduces *operation* over a contiguous subsequence of items in the array.

Parameters

- **capacity** – (int) Total size of the array - must be a power of two.
- **operation** – (lambda obj, obj -> obj) and operation for combining elements (eg. sum, max) must form a mathematical group together with the set of possible values for array elements (i.e. be associative)
- **neutral_element** – (obj) neutral element for the operation above. eg. float('-inf') for max and 0 for sum.

`__module__ = 'rlzoo.common.buffer'`

`find_prefixsum_idx(prefixsum)`

Find the highest index i in the array such that $\text{sum}(\text{arr}[0] + \text{arr}[1] + \dots + \text{arr}[i - 1]) \leq \text{prefixsum}$

if array values are probabilities, this function allows to sample indexes according to the discrete probability efficiently.

Parameters **prefixsum** – (float) upperbound on the sum of array prefix

Returns:

idx: (int) highest index satisfying the prefixsum constraint

`sum(start=0, end=None)`

Returns $\text{arr}[\text{start}] + \dots + \text{arr}[\text{end}]$

19.1 Distributions for Stochastic Policy in RLzoo

Definition of parametrized distributions. Adapted from openai/baselines

```
class rlzoo.common.distributions.Categorical (ndim, logits=None)
    Bases: rlzoo.common.distributions.Distribution

    Creates a categorical distribution

    entropy ()
        Calculate the entropy of distribution.

    get_param ()

    greedy_sample ()
        Get actions greedily

    kl (logits)
        Args: logits (tensor): logits variables of another distribution

    logp (x)
        Calculate log probability of a sample.

    ndim

    neglogp (x)
        Calculate negative log probability of a sample.

    sample ()
        Sample actions from distribution, using the Gumbel-Softmax trick

    set_param (logits)
        Args: logits (tensor): logits variables to set

class rlzoo.common.distributions.DiagGaussian (ndim, mean_logstd=None)
    Bases: rlzoo.common.distributions.Distribution
```

Creates a diagonal Gaussian distribution

entropy ()

Calculate the entropy of distribution.

get_param ()

Get parameters

greedy_sample ()

Get actions greedily/deterministically

kl (*mean_logstd*)

Args: *mean_logstd* (tensor): mean and logstd of another distribution

logp (*x*)

Calculate log probability of a sample.

ndim

neglogp (*x*)

Calculate negative log probability of a sample.

sample ()

Get actions in deterministic or stochastic manner

set_param (*mean_logstd*)

Args: *mean_logstd* (tensor): mean and log std

class rlzoo.common.distributions.**Distribution**

Bases: object

A particular probability distribution

entropy ()

Calculate the entropy of distribution.

kl (**parameters*)

Calculate Kullback–Leibler divergence

logp (*x*)

Calculate log probability of a sample.

neglogp (*x*)

Calculate negative log probability of a sample.

sample (**args, **kwargs*)

Sampling from distribution. Allow explore parameters.

set_param (**args, **kwargs*)

rlzoo.common.distributions.**expand_dims** (*func*)

rlzoo.common.distributions.**make_dist** (*ac_space*)

Get distribution based on action space

Parameters *ac_space* – gym.spaces.Space

20.1 Environment Wrappers in RLzoo

Env wrappers Most common wrappers can be checked from following links for usage:

<https://pypi.org/project/gym-vec-env>

<https://github.com/openai/baselines/blob/master/baselines/common/wrappers.py>

```
rlzoo.common.env_wrappers.build_env(env_id, env_type, vectorized=False, seed=0, re-  
ward_shaping=None, nenv=1, **kwargs)
```

Build env based on options

Parameters

- **env_id** – (str) environment id
- **env_type** – (str) atari, classic_control, box2d
- **vectorized** – (bool) whether sampling parallel
- **seed** – (int) random seed for env
- **reward_shaping** – (callable) callable function for reward shaping
- **nenv** – (int) how many processes will be used in sampling
- **kwargs** – (dict)
- **max_episode_steps** – (int) the maximum episode steps

```
class rlzoo.common.env_wrappers.TimeLimit(env, max_episode_steps=None)
```

Bases: gym.core.Wrapper

```
reset (**kwargs)
```

Resets the state of the environment and returns an initial observation.

Returns: observation (object): the initial observation.

step (*ac*)

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).

Args: action (object): an action provided by the agent

Returns: observation (object): 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)

```
class rlzoo.common.env_wrappers.NoopResetEnv (env, noop_max=30)
```

Bases: gym.core.Wrapper

reset (***kwargs*)

Do no-op action for a number of steps in [1, noop_max].

step (*ac*)

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).

Args: action (object): an action provided by the agent

Returns: observation (object): 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)

```
class rlzoo.common.env_wrappers.FireResetEnv (env)
```

Bases: gym.core.Wrapper

reset (***kwargs*)

Resets the state of the environment and returns an initial observation.

Returns: observation (object): the initial observation.

step (*ac*)

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).

Args: action (object): an action provided by the agent

Returns: observation (object): 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)

```
class rlzoo.common.env_wrappers.EpisodicLifeEnv (env)
```

Bases: gym.core.Wrapper

reset (***kwargs*)

Reset only when lives are exhausted. This way all states are still reachable even though lives are episodic, and the learner need not know about any of this behind-the-scenes.

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).

Args: action (object): an action provided by the agent

Returns: observation (object): 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)

```
class rlzoo.common.env_wrappers.MaxAndSkipEnv (env, skip=4)
    Bases: gym.core.Wrapper
```

```
reset ( **kwargs)
    Resets the state of the environment and returns an initial observation.
```

Returns: observation (object): the initial observation.

```
step (action)
    Repeat action, sum reward, and max over last observations.
```

```
class rlzoo.common.env_wrappers.ClipRewardEnv (env)
    Bases: gym.core.RewardWrapper
```

```
reward (reward)
    Bin reward to {+1, 0, -1} by its sign.
```

```
class rlzoo.common.env_wrappers.WarpFrame (env, width=84, height=84, grayscale=True)
    Bases: gym.core.ObservationWrapper
```

```
observation (frame)
```

```
class rlzoo.common.env_wrappers.FrameStack (env, k)
    Bases: gym.core.Wrapper
```

```
reset ()
    Resets the state of the environment and returns an initial observation.
```

Returns: observation (object): the initial observation.

```
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).

Args: action (object): an action provided by the agent

Returns: observation (object): 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)

```
class rlzoo.common.env_wrappers.LazyFrames (frames)
    Bases: object
```

```
class rlzoo.common.env_wrappers.RewardShaping (env, func)
    Bases: gym.core.RewardWrapper
```

Shaping the reward For reward scale, func can be *lambda r: r * scale*

```
reward (reward)
```

```
class rlzoo.common.env_wrappers.SubprocVecEnv (env_fns)
    Bases: object
```

```
close ()

reset ()
    Reset all the environments and return an array of observations, or a tuple of observation arrays. If
    step_async is still doing work, that work will be cancelled and step_wait() should not be called until
    step_async() is invoked again.

step (actions)

class rlzoo.common.env_wrappers.VecFrameStack (env, k)
    Bases: object

    reset ()

    step (action)

class rlzoo.common.env_wrappers.Monitor (env, info_keywords=None)
    Bases: gym.core.Wrapper

    reset (**kwargs)
        Resets the state of the environment and returns an initial observation.

        Returns: observation (object): the initial observation.

    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).

        Args: action (object): an action provided by the agent

        Returns: observation (object): 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)

class rlzoo.common.env_wrappers.NormalizedActions (env)
    Bases: gym.core.ActionWrapper

class rlzoo.common.env_wrappers.DmObsTrans (env)
    Bases: gym.core.Wrapper

    Observation process for DeepMind Control Suite environments

    reset (**kwargs)
        Resets the state of the environment and returns an initial observation.

        Returns: observation (object): the initial observation.

    step (ac)
        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).

        Args: action (object): an action provided by the agent

        Returns: observation (object): 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)
```

21.1 List of Supported Environments in RLzoo

`rlzoo.common.env_list.get_envlist(env_type)`
get list of env names wrt the type of env

```
1 all_env_list = {  
2     ## Gym  
3     # Atari  
4     'atari': ['AirRaid-v0',  
5               'AirRaid-v4',  
6               'AirRaidDeterministic-v0',  
7               'AirRaidDeterministic-v4',  
8               'AirRaidNoFrameskip-v0',  
9               'AirRaidNoFrameskip-v4',  
10              'AirRaid-ram-v0',  
11              'AirRaid-ram-v4',  
12              'AirRaid-ramDeterministic-v0',  
13              'AirRaid-ramDeterministic-v4',  
14              'AirRaid-ramNoFrameskip-v0',  
15              'AirRaid-ramNoFrameskip-v4',  
16              'Alien-v0',  
17              'Alien-v4',  
18              'AlienDeterministic-v0',  
19              'AlienDeterministic-v4',  
20              'AlienNoFrameskip-v0',  
21              'AlienNoFrameskip-v4',  
22              'Alien-ram-v0',  
23              'Alien-ram-v4',  
24              'Alien-ramDeterministic-v0',  
25              'Alien-ramDeterministic-v4',  
26              'Alien-ramNoFrameskip-v0',  
27              'Alien-ramNoFrameskip-v4',  
28              'Amidar-v0',
```

(continues on next page)

(continued from previous page)

```

29     'Amidar-v4',
30     'AmidarDeterministic-v0',
31     'AmidarDeterministic-v4',
32     'AmidarNoFrameskip-v0',
33     'AmidarNoFrameskip-v4',
34     'Amidar-ram-v0',
35     'Amidar-ram-v4',
36     'Amidar-ramDeterministic-v0',
37     'Amidar-ramDeterministic-v4',
38     'Amidar-ramNoFrameskip-v0',
39     'Amidar-ramNoFrameskip-v4',
40     'Assault-v0',
41     'Assault-v4',
42     'AssaultDeterministic-v0',
43     'AssaultDeterministic-v4',
44     'AssaultNoFrameskip-v0',
45     'AssaultNoFrameskip-v4',
46     'Assault-ram-v0',
47     'Assault-ram-v4',
48     'Assault-ramDeterministic-v0',
49     'Assault-ramDeterministic-v4',
50     'Assault-ramNoFrameskip-v0',
51     'Assault-ramNoFrameskip-v4',
52     'Asterix-v0',
53     'Asterix-v4',
54     'AsterixDeterministic-v0',
55     'AsterixDeterministic-v4',
56     'AsterixNoFrameskip-v0',
57     'AsterixNoFrameskip-v4',
58     'Asterix-ram-v0',
59     'Asterix-ram-v4',
60     'Asterix-ramDeterministic-v0',
61     'Asterix-ramDeterministic-v4',
62     'Asterix-ramNoFrameskip-v0',
63     'Asterix-ramNoFrameskip-v4',
64     'Asteroids-v0',
65     'Asteroids-v4',
66     'AsteroidsDeterministic-v0',
67     'AsteroidsDeterministic-v4',
68     'AsteroidsNoFrameskip-v0',
69     'AsteroidsNoFrameskip-v4',
70     'Asteroids-ram-v0',
71     'Asteroids-ram-v4',
72     'Asteroids-ramDeterministic-v0',
73     'Asteroids-ramDeterministic-v4',
74     'Asteroids-ramNoFrameskip-v0',
75     'Asteroids-ramNoFrameskip-v4',
76     'Atlantis-v0',
77     'Atlantis-v4',
78     'AtlantisDeterministic-v0',
79     'AtlantisDeterministic-v4',
80     'AtlantisNoFrameskip-v0',
81     'AtlantisNoFrameskip-v4',
82     'Atlantis-ram-v0',
83     'Atlantis-ram-v4',
84     'Atlantis-ramDeterministic-v0',
85     'Atlantis-ramDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

86      'Atlantis-ramNoFrameskip-v0',
87      'Atlantis-ramNoFrameskip-v4',
88      'BankHeist-v0',
89      'BankHeist-v4',
90      'BankHeistDeterministic-v0',
91      'BankHeistDeterministic-v4',
92      'BankHeistNoFrameskip-v0',
93      'BankHeistNoFrameskip-v4',
94      'BankHeist-ram-v0',
95      'BankHeist-ram-v4',
96      'BankHeist-ramDeterministic-v0',
97      'BankHeist-ramDeterministic-v4',
98      'BankHeist-ramNoFrameskip-v0',
99      'BankHeist-ramNoFrameskip-v4',
100     'BattleZone-v0',
101     'BattleZone-v4',
102     'BattleZoneDeterministic-v0',
103     'BattleZoneDeterministic-v4',
104     'BattleZoneNoFrameskip-v0',
105     'BattleZoneNoFrameskip-v4',
106     'BattleZone-ram-v0',
107     'BattleZone-ram-v4',
108     'BattleZone-ramDeterministic-v0',
109     'BattleZone-ramDeterministic-v4',
110     'BattleZone-ramNoFrameskip-v0',
111     'BattleZone-ramNoFrameskip-v4',
112     'BeamRider-v0',
113     'BeamRider-v4',
114     'BeamRiderDeterministic-v0',
115     'BeamRiderDeterministic-v4',
116     'BeamRiderNoFrameskip-v0',
117     'BeamRiderNoFrameskip-v4',
118     'BeamRider-ram-v0',
119     'BeamRider-ram-v4',
120     'BeamRider-ramDeterministic-v0',
121     'BeamRider-ramDeterministic-v4',
122     'BeamRider-ramNoFrameskip-v0',
123     'BeamRider-ramNoFrameskip-v4',
124     'Berzerk-v0',
125     'Berzerk-v4',
126     'BerzerkDeterministic-v0',
127     'BerzerkDeterministic-v4',
128     'BerzerkNoFrameskip-v0',
129     'BerzerkNoFrameskip-v4',
130     'Berzerk-ram-v0',
131     'Berzerk-ram-v4',
132     'Berzerk-ramDeterministic-v0',
133     'Berzerk-ramDeterministic-v4',
134     'Berzerk-ramNoFrameskip-v0',
135     'Berzerk-ramNoFrameskip-v4',
136     'Bowling-v0',
137     'Bowling-v4',
138     'BowlingDeterministic-v0',
139     'BowlingDeterministic-v4',
140     'BowlingNoFrameskip-v0',
141     'BowlingNoFrameskip-v4',
142     'Bowling-ram-v0',

```

(continues on next page)

(continued from previous page)

```

143     'Bowling-ram-v4',
144     'Bowling-ramDeterministic-v0',
145     'Bowling-ramDeterministic-v4',
146     'Bowling-ramNoFrameskip-v0',
147     'Bowling-ramNoFrameskip-v4',
148     'Boxing-v0',
149     'Boxing-v4',
150     'BoxingDeterministic-v0',
151     'BoxingDeterministic-v4',
152     'BoxingNoFrameskip-v0',
153     'BoxingNoFrameskip-v4',
154     'Boxing-ram-v0',
155     'Boxing-ram-v4',
156     'Boxing-ramDeterministic-v0',
157     'Boxing-ramDeterministic-v4',
158     'Boxing-ramNoFrameskip-v0',
159     'Boxing-ramNoFrameskip-v4',
160     'Breakout-v0',
161     'Breakout-v4',
162     'BreakoutDeterministic-v0',
163     'BreakoutDeterministic-v4',
164     'BreakoutNoFrameskip-v0',
165     'BreakoutNoFrameskip-v4',
166     'Breakout-ram-v0',
167     'Breakout-ram-v4',
168     'Breakout-ramDeterministic-v0',
169     'Breakout-ramDeterministic-v4',
170     'Breakout-ramNoFrameskip-v0',
171     'Breakout-ramNoFrameskip-v4',
172     'Carnival-v0',
173     'Carnival-v4',
174     'CarnivalDeterministic-v0',
175     'CarnivalDeterministic-v4',
176     'CarnivalNoFrameskip-v0',
177     'CarnivalNoFrameskip-v4',
178     'Carnival-ram-v0',
179     'Carnival-ram-v4',
180     'Carnival-ramDeterministic-v0',
181     'Carnival-ramDeterministic-v4',
182     'Carnival-ramNoFrameskip-v0',
183     'Carnival-ramNoFrameskip-v4',
184     'Centipede-v0',
185     'Centipede-v4',
186     'CentipedeDeterministic-v0',
187     'CentipedeDeterministic-v4',
188     'CentipedeNoFrameskip-v0',
189     'CentipedeNoFrameskip-v4',
190     'Centipede-ram-v0',
191     'Centipede-ram-v4',
192     'Centipede-ramDeterministic-v0',
193     'Centipede-ramDeterministic-v4',
194     'Centipede-ramNoFrameskip-v0',
195     'Centipede-ramNoFrameskip-v4',
196     'ChopperCommand-v0',
197     'ChopperCommand-v4',
198     'ChopperCommandDeterministic-v0',
199     'ChopperCommandDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

200     'ChopperCommandNoFrameskip-v0',
201     'ChopperCommandNoFrameskip-v4',
202     'ChopperCommand-ram-v0',
203     'ChopperCommand-ram-v4',
204     'ChopperCommand-ramDeterministic-v0',
205     'ChopperCommand-ramDeterministic-v4',
206     'ChopperCommand-ramNoFrameskip-v0',
207     'ChopperCommand-ramNoFrameskip-v4',
208     'CrazyClimber-v0',
209     'CrazyClimber-v4',
210     'CrazyClimberDeterministic-v0',
211     'CrazyClimberDeterministic-v4',
212     'CrazyClimberNoFrameskip-v0',
213     'CrazyClimberNoFrameskip-v4',
214     'CrazyClimber-ram-v0',
215     'CrazyClimber-ram-v4',
216     'CrazyClimber-ramDeterministic-v0',
217     'CrazyClimber-ramDeterministic-v4',
218     'CrazyClimber-ramNoFrameskip-v0',
219     'CrazyClimber-ramNoFrameskip-v4',
220     'DemonAttack-v0',
221     'DemonAttack-v4',
222     'DemonAttackDeterministic-v0',
223     'DemonAttackDeterministic-v4',
224     'DemonAttackNoFrameskip-v0',
225     'DemonAttackNoFrameskip-v4',
226     'DemonAttack-ram-v0',
227     'DemonAttack-ram-v4',
228     'DemonAttack-ramDeterministic-v0',
229     'DemonAttack-ramDeterministic-v4',
230     'DemonAttack-ramNoFrameskip-v0',
231     'DemonAttack-ramNoFrameskip-v4',
232     'DoubleDunk-v0',
233     'DoubleDunk-v4',
234     'DoubleDunkDeterministic-v0',
235     'DoubleDunkDeterministic-v4',
236     'DoubleDunkNoFrameskip-v0',
237     'DoubleDunkNoFrameskip-v4',
238     'DoubleDunk-ram-v0',
239     'DoubleDunk-ram-v4',
240     'DoubleDunk-ramDeterministic-v0',
241     'DoubleDunk-ramDeterministic-v4',
242     'DoubleDunk-ramNoFrameskip-v0',
243     'DoubleDunk-ramNoFrameskip-v4',
244     'ElevatorAction-v0',
245     'ElevatorAction-v4',
246     'ElevatorActionDeterministic-v0',
247     'ElevatorActionDeterministic-v4',
248     'ElevatorActionNoFrameskip-v0',
249     'ElevatorActionNoFrameskip-v4',
250     'ElevatorAction-ram-v0',
251     'ElevatorAction-ram-v4',
252     'ElevatorAction-ramDeterministic-v0',
253     'ElevatorAction-ramDeterministic-v4',
254     'ElevatorAction-ramNoFrameskip-v0',
255     'ElevatorAction-ramNoFrameskip-v4',
256     'Enduro-v0',

```

(continues on next page)

(continued from previous page)

```

257         'Enduro-v4',
258         'EnduroDeterministic-v0',
259         'EnduroDeterministic-v4',
260         'EnduroNoFrameskip-v0',
261         'EnduroNoFrameskip-v4',
262         'Enduro-ram-v0',
263         'Enduro-ram-v4',
264         'Enduro-ramDeterministic-v0',
265         'Enduro-ramDeterministic-v4',
266         'Enduro-ramNoFrameskip-v0',
267         'Enduro-ramNoFrameskip-v4',
268         'FishingDerby-v0',
269         'FishingDerby-v4',
270         'FishingDerbyDeterministic-v0',
271         'FishingDerbyDeterministic-v4',
272         'FishingDerbyNoFrameskip-v0',
273         'FishingDerbyNoFrameskip-v4',
274         'FishingDerby-ram-v0',
275         'FishingDerby-ram-v4',
276         'FishingDerby-ramDeterministic-v0',
277         'FishingDerby-ramDeterministic-v4',
278         'FishingDerby-ramNoFrameskip-v0',
279         'FishingDerby-ramNoFrameskip-v4',
280         'Freeway-v0',
281         'Freeway-v4',
282         'FreewayDeterministic-v0',
283         'FreewayDeterministic-v4',
284         'FreewayNoFrameskip-v0',
285         'FreewayNoFrameskip-v4',
286         'Freeway-ram-v0',
287         'Freeway-ram-v4',
288         'Freeway-ramDeterministic-v0',
289         'Freeway-ramDeterministic-v4',
290         'Freeway-ramNoFrameskip-v0',
291         'Freeway-ramNoFrameskip-v4',
292         'Frostbite-v0',
293         'Frostbite-v4',
294         'FrostbiteDeterministic-v0',
295         'FrostbiteDeterministic-v4',
296         'FrostbiteNoFrameskip-v0',
297         'FrostbiteNoFrameskip-v4',
298         'Frostbite-ram-v0',
299         'Frostbite-ram-v4',
300         'Frostbite-ramDeterministic-v0',
301         'Frostbite-ramDeterministic-v4',
302         'Frostbite-ramNoFrameskip-v0',
303         'Frostbite-ramNoFrameskip-v4',
304         'Gopher-v0',
305         'Gopher-v4',
306         'GopherDeterministic-v0',
307         'GopherDeterministic-v4',
308         'GopherNoFrameskip-v0',
309         'GopherNoFrameskip-v4',
310         'Gopher-ram-v0',
311         'Gopher-ram-v4',
312         'Gopher-ramDeterministic-v0',
313         'Gopher-ramDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

314     'Gopher-ramNoFrameskip-v0',
315     'Gopher-ramNoFrameskip-v4',
316     'Gravitar-v0',
317     'Gravitar-v4',
318     'GravitarDeterministic-v0',
319     'GravitarDeterministic-v4',
320     'GravitarNoFrameskip-v0',
321     'GravitarNoFrameskip-v4',
322     'Gravitar-ram-v0',
323     'Gravitar-ram-v4',
324     'Gravitar-ramDeterministic-v0',
325     'Gravitar-ramDeterministic-v4',
326     'Gravitar-ramNoFrameskip-v0',
327     'Gravitar-ramNoFrameskip-v4',
328     'Hero-v0',
329     'Hero-v4',
330     'HeroDeterministic-v0',
331     'HeroDeterministic-v4',
332     'HeroNoFrameskip-v0',
333     'HeroNoFrameskip-v4',
334     'Hero-ram-v0',
335     'Hero-ram-v4',
336     'Hero-ramDeterministic-v0',
337     'Hero-ramDeterministic-v4',
338     'Hero-ramNoFrameskip-v0',
339     'Hero-ramNoFrameskip-v4',
340     'IceHockey-v0',
341     'IceHockey-v4',
342     'IceHockeyDeterministic-v0',
343     'IceHockeyDeterministic-v4',
344     'IceHockeyNoFrameskip-v0',
345     'IceHockeyNoFrameskip-v4',
346     'IceHockey-ram-v0',
347     'IceHockey-ram-v4',
348     'IceHockey-ramDeterministic-v0',
349     'IceHockey-ramDeterministic-v4',
350     'IceHockey-ramNoFrameskip-v0',
351     'IceHockey-ramNoFrameskip-v4',
352     'Jamesbond-v0',
353     'Jamesbond-v4',
354     'JamesbondDeterministic-v0',
355     'JamesbondDeterministic-v4',
356     'JamesbondNoFrameskip-v0',
357     'JamesbondNoFrameskip-v4',
358     'Jamesbond-ram-v0',
359     'Jamesbond-ram-v4',
360     'Jamesbond-ramDeterministic-v0',
361     'Jamesbond-ramDeterministic-v4',
362     'Jamesbond-ramNoFrameskip-v0',
363     'Jamesbond-ramNoFrameskip-v4',
364     'JourneyEscape-v0',
365     'JourneyEscape-v4',
366     'JourneyEscapeDeterministic-v0',
367     'JourneyEscapeDeterministic-v4',
368     'JourneyEscapeNoFrameskip-v0',
369     'JourneyEscapeNoFrameskip-v4',
370     'JourneyEscape-ram-v0',

```

(continues on next page)

(continued from previous page)

```

371     'JourneyEscape-ram-v4',
372     'JourneyEscape-ramDeterministic-v0',
373     'JourneyEscape-ramDeterministic-v4',
374     'JourneyEscape-ramNoFrameskip-v0',
375     'JourneyEscape-ramNoFrameskip-v4',
376     'Kangaroo-v0',
377     'Kangaroo-v4',
378     'KangarooDeterministic-v0',
379     'KangarooDeterministic-v4',
380     'KangarooNoFrameskip-v0',
381     'KangarooNoFrameskip-v4',
382     'Kangaroo-ram-v0',
383     'Kangaroo-ram-v4',
384     'Kangaroo-ramDeterministic-v0',
385     'Kangaroo-ramDeterministic-v4',
386     'Kangaroo-ramNoFrameskip-v0',
387     'Kangaroo-ramNoFrameskip-v4',
388     'Krull-v0',
389     'Krull-v4',
390     'KrullDeterministic-v0',
391     'KrullDeterministic-v4',
392     'KrullNoFrameskip-v0',
393     'KrullNoFrameskip-v4',
394     'Krull-ram-v0',
395     'Krull-ram-v4',
396     'Krull-ramDeterministic-v0',
397     'Krull-ramDeterministic-v4',
398     'Krull-ramNoFrameskip-v0',
399     'Krull-ramNoFrameskip-v4',
400     'KungFuMaster-v0',
401     'KungFuMaster-v4',
402     'KungFuMasterDeterministic-v0',
403     'KungFuMasterDeterministic-v4',
404     'KungFuMasterNoFrameskip-v0',
405     'KungFuMasterNoFrameskip-v4',
406     'KungFuMaster-ram-v0',
407     'KungFuMaster-ram-v4',
408     'KungFuMaster-ramDeterministic-v0',
409     'KungFuMaster-ramDeterministic-v4',
410     'KungFuMaster-ramNoFrameskip-v0',
411     'KungFuMaster-ramNoFrameskip-v4',
412     'MontezumaRevenge-v0',
413     'MontezumaRevenge-v4',
414     'MontezumaRevengeDeterministic-v0',
415     'MontezumaRevengeDeterministic-v4',
416     'MontezumaRevengeNoFrameskip-v0',
417     'MontezumaRevengeNoFrameskip-v4',
418     'MontezumaRevenge-ram-v0',
419     'MontezumaRevenge-ram-v4',
420     'MontezumaRevenge-ramDeterministic-v0',
421     'MontezumaRevenge-ramDeterministic-v4',
422     'MontezumaRevenge-ramNoFrameskip-v0',
423     'MontezumaRevenge-ramNoFrameskip-v4',
424     'MsPacman-v0',
425     'MsPacman-v4',
426     'MsPacmanDeterministic-v0',
427     'MsPacmanDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

428     'MsPacmanNoFrameskip-v0',
429     'MsPacmanNoFrameskip-v4',
430     'MsPacman-ram-v0',
431     'MsPacman-ram-v4',
432     'MsPacman-ramDeterministic-v0',
433     'MsPacman-ramDeterministic-v4',
434     'MsPacman-ramNoFrameskip-v0',
435     'MsPacman-ramNoFrameskip-v4',
436     'NameThisGame-v0',
437     'NameThisGame-v4',
438     'NameThisGameDeterministic-v0',
439     'NameThisGameDeterministic-v4',
440     'NameThisGameNoFrameskip-v0',
441     'NameThisGameNoFrameskip-v4',
442     'NameThisGame-ram-v0',
443     'NameThisGame-ram-v4',
444     'NameThisGame-ramDeterministic-v0',
445     'NameThisGame-ramDeterministic-v4',
446     'NameThisGame-ramNoFrameskip-v0',
447     'NameThisGame-ramNoFrameskip-v4',
448     'Phoenix-v0',
449     'Phoenix-v4',
450     'PhoenixDeterministic-v0',
451     'PhoenixDeterministic-v4',
452     'PhoenixNoFrameskip-v0',
453     'PhoenixNoFrameskip-v4',
454     'Phoenix-ram-v0',
455     'Phoenix-ram-v4',
456     'Phoenix-ramDeterministic-v0',
457     'Phoenix-ramDeterministic-v4',
458     'Phoenix-ramNoFrameskip-v0',
459     'Phoenix-ramNoFrameskip-v4',
460     'Pitfall-v0',
461     'Pitfall-v4',
462     'PitfallDeterministic-v0',
463     'PitfallDeterministic-v4',
464     'PitfallNoFrameskip-v0',
465     'PitfallNoFrameskip-v4',
466     'Pitfall-ram-v0',
467     'Pitfall-ram-v4',
468     'Pitfall-ramDeterministic-v0',
469     'Pitfall-ramDeterministic-v4',
470     'Pitfall-ramNoFrameskip-v0',
471     'Pitfall-ramNoFrameskip-v4',
472     'Pong-v0',
473     'Pong-v4',
474     'PongDeterministic-v0',
475     'PongDeterministic-v4',
476     'PongNoFrameskip-v0',
477     'PongNoFrameskip-v4',
478     'Pong-ram-v0',
479     'Pong-ram-v4',
480     'Pong-ramDeterministic-v0',
481     'Pong-ramDeterministic-v4',
482     'Pong-ramNoFrameskip-v0',
483     'Pong-ramNoFrameskip-v4',
484     'Pooyan-v0',

```

(continues on next page)

(continued from previous page)

```

485     'Pooyan-v4',
486     'PooyanDeterministic-v0',
487     'PooyanDeterministic-v4',
488     'PooyanNoFrameskip-v0',
489     'PooyanNoFrameskip-v4',
490     'Pooyan-ram-v0',
491     'Pooyan-ram-v4',
492     'Pooyan-ramDeterministic-v0',
493     'Pooyan-ramDeterministic-v4',
494     'Pooyan-ramNoFrameskip-v0',
495     'Pooyan-ramNoFrameskip-v4',
496     'PrivateEye-v0',
497     'PrivateEye-v4',
498     'PrivateEyeDeterministic-v0',
499     'PrivateEyeDeterministic-v4',
500     'PrivateEyeNoFrameskip-v0',
501     'PrivateEyeNoFrameskip-v4',
502     'PrivateEye-ram-v0',
503     'PrivateEye-ram-v4',
504     'PrivateEye-ramDeterministic-v0',
505     'PrivateEye-ramDeterministic-v4',
506     'PrivateEye-ramNoFrameskip-v0',
507     'PrivateEye-ramNoFrameskip-v4',
508     'Qbert-v0',
509     'Qbert-v4',
510     'QbertDeterministic-v0',
511     'QbertDeterministic-v4',
512     'QbertNoFrameskip-v0',
513     'QbertNoFrameskip-v4',
514     'Qbert-ram-v0',
515     'Qbert-ram-v4',
516     'Qbert-ramDeterministic-v0',
517     'Qbert-ramDeterministic-v4',
518     'Qbert-ramNoFrameskip-v0',
519     'Qbert-ramNoFrameskip-v4',
520     'Riverraid-v0',
521     'Riverraid-v4',
522     'RiverraidDeterministic-v0',
523     'RiverraidDeterministic-v4',
524     'RiverraidNoFrameskip-v0',
525     'RiverraidNoFrameskip-v4',
526     'Riverraid-ram-v0',
527     'Riverraid-ram-v4',
528     'Riverraid-ramDeterministic-v0',
529     'Riverraid-ramDeterministic-v4',
530     'Riverraid-ramNoFrameskip-v0',
531     'Riverraid-ramNoFrameskip-v4',
532     'RoadRunner-v0',
533     'RoadRunner-v4',
534     'RoadRunnerDeterministic-v0',
535     'RoadRunnerDeterministic-v4',
536     'RoadRunnerNoFrameskip-v0',
537     'RoadRunnerNoFrameskip-v4',
538     'RoadRunner-ram-v0',
539     'RoadRunner-ram-v4',
540     'RoadRunner-ramDeterministic-v0',
541     'RoadRunner-ramDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

542     'RoadRunner-ramNoFrameskip-v0',
543     'RoadRunner-ramNoFrameskip-v4',
544     'Robotank-v0',
545     'Robotank-v4',
546     'RobotankDeterministic-v0',
547     'RobotankDeterministic-v4',
548     'RobotankNoFrameskip-v0',
549     'RobotankNoFrameskip-v4',
550     'Robotank-ram-v0',
551     'Robotank-ram-v4',
552     'Robotank-ramDeterministic-v0',
553     'Robotank-ramDeterministic-v4',
554     'Robotank-ramNoFrameskip-v0',
555     'Robotank-ramNoFrameskip-v4',
556     'Seaquest-v0',
557     'Seaquest-v4',
558     'SeaquestDeterministic-v0',
559     'SeaquestDeterministic-v4',
560     'SeaquestNoFrameskip-v0',
561     'SeaquestNoFrameskip-v4',
562     'Seaquest-ram-v0',
563     'Seaquest-ram-v4',
564     'Seaquest-ramDeterministic-v0',
565     'Seaquest-ramDeterministic-v4',
566     'Seaquest-ramNoFrameskip-v0',
567     'Seaquest-ramNoFrameskip-v4',
568     'Skiing-v0',
569     'Skiing-v4',
570     'SkiingDeterministic-v0',
571     'SkiingDeterministic-v4',
572     'SkiingNoFrameskip-v0',
573     'SkiingNoFrameskip-v4',
574     'Skiing-ram-v0',
575     'Skiing-ram-v4',
576     'Skiing-ramDeterministic-v0',
577     'Skiing-ramDeterministic-v4',
578     'Skiing-ramNoFrameskip-v0',
579     'Skiing-ramNoFrameskip-v4',
580     'Solaris-v0',
581     'Solaris-v4',
582     'SolarisDeterministic-v0',
583     'SolarisDeterministic-v4',
584     'SolarisNoFrameskip-v0',
585     'SolarisNoFrameskip-v4',
586     'Solaris-ram-v0',
587     'Solaris-ram-v4',
588     'Solaris-ramDeterministic-v0',
589     'Solaris-ramDeterministic-v4',
590     'Solaris-ramNoFrameskip-v0',
591     'Solaris-ramNoFrameskip-v4',
592     'SpaceInvaders-v0',
593     'SpaceInvaders-v4',
594     'SpaceInvadersDeterministic-v0',
595     'SpaceInvadersDeterministic-v4',
596     'SpaceInvadersNoFrameskip-v0',
597     'SpaceInvadersNoFrameskip-v4',
598     'SpaceInvaders-ram-v0',

```

(continues on next page)

(continued from previous page)

```

599     'SpaceInvaders-ram-v4',
600     'SpaceInvaders-ramDeterministic-v0',
601     'SpaceInvaders-ramDeterministic-v4',
602     'SpaceInvaders-ramNoFrameskip-v0',
603     'SpaceInvaders-ramNoFrameskip-v4',
604     'StarGunner-v0',
605     'StarGunner-v4',
606     'StarGunnerDeterministic-v0',
607     'StarGunnerDeterministic-v4',
608     'StarGunnerNoFrameskip-v0',
609     'StarGunnerNoFrameskip-v4',
610     'StarGunner-ram-v0',
611     'StarGunner-ram-v4',
612     'StarGunner-ramDeterministic-v0',
613     'StarGunner-ramDeterministic-v4',
614     'StarGunner-ramNoFrameskip-v0',
615     'StarGunner-ramNoFrameskip-v4',
616     'Tennis-v0',
617     'Tennis-v4',
618     'TennisDeterministic-v0',
619     'TennisDeterministic-v4',
620     'TennisNoFrameskip-v0',
621     'TennisNoFrameskip-v4',
622     'Tennis-ram-v0',
623     'Tennis-ram-v4',
624     'Tennis-ramDeterministic-v0',
625     'Tennis-ramDeterministic-v4',
626     'Tennis-ramNoFrameskip-v0',
627     'Tennis-ramNoFrameskip-v4',
628     'TimePilot-v0',
629     'TimePilot-v4',
630     'TimePilotDeterministic-v0',
631     'TimePilotDeterministic-v4',
632     'TimePilotNoFrameskip-v0',
633     'TimePilotNoFrameskip-v4',
634     'TimePilot-ram-v0',
635     'TimePilot-ram-v4',
636     'TimePilot-ramDeterministic-v0',
637     'TimePilot-ramDeterministic-v4',
638     'TimePilot-ramNoFrameskip-v0',
639     'TimePilot-ramNoFrameskip-v4',
640     'Tutankham-v0',
641     'Tutankham-v4',
642     'TutankhamDeterministic-v0',
643     'TutankhamDeterministic-v4',
644     'TutankhamNoFrameskip-v0',
645     'TutankhamNoFrameskip-v4',
646     'Tutankham-ram-v0',
647     'Tutankham-ram-v4',
648     'Tutankham-ramDeterministic-v0',
649     'Tutankham-ramDeterministic-v4',
650     'Tutankham-ramNoFrameskip-v0',
651     'Tutankham-ramNoFrameskip-v4',
652     'UpNDown-v0',
653     'UpNDown-v4',
654     'UpNDownDeterministic-v0',
655     'UpNDownDeterministic-v4',

```

(continues on next page)

(continued from previous page)

```

656     'UpNDownNoFrameskip-v0',
657     'UpNDownNoFrameskip-v4',
658     'UpNDown-ram-v0',
659     'UpNDown-ram-v4',
660     'UpNDown-ramDeterministic-v0',
661     'UpNDown-ramDeterministic-v4',
662     'UpNDown-ramNoFrameskip-v0',
663     'UpNDown-ramNoFrameskip-v4',
664     'Venture-v0',
665     'Venture-v4',
666     'VentureDeterministic-v0',
667     'VentureDeterministic-v4',
668     'VentureNoFrameskip-v0',
669     'VentureNoFrameskip-v4',
670     'Venture-ram-v0',
671     'Venture-ram-v4',
672     'Venture-ramDeterministic-v0',
673     'Venture-ramDeterministic-v4',
674     'Venture-ramNoFrameskip-v0',
675     'Venture-ramNoFrameskip-v4',
676     'VideoPinball-v0',
677     'VideoPinball-v4',
678     'VideoPinballDeterministic-v0',
679     'VideoPinballDeterministic-v4',
680     'VideoPinballNoFrameskip-v0',
681     'VideoPinballNoFrameskip-v4',
682     'VideoPinball-ram-v0',
683     'VideoPinball-ram-v4',
684     'VideoPinball-ramDeterministic-v0',
685     'VideoPinball-ramDeterministic-v4',
686     'VideoPinball-ramNoFrameskip-v0',
687     'VideoPinball-ramNoFrameskip-v4',
688     'WizardOfWor-v0',
689     'WizardOfWor-v4',
690     'WizardOfWorDeterministic-v0',
691     'WizardOfWorDeterministic-v4',
692     'WizardOfWorNoFrameskip-v0',
693     'WizardOfWorNoFrameskip-v4',
694     'WizardOfWor-ram-v0',
695     'WizardOfWor-ram-v4',
696     'WizardOfWor-ramDeterministic-v0',
697     'WizardOfWor-ramDeterministic-v4',
698     'WizardOfWor-ramNoFrameskip-v0',
699     'WizardOfWor-ramNoFrameskip-v4',
700     'YarsRevenge-v0',
701     'YarsRevenge-v4',
702     'YarsRevengeDeterministic-v0',
703     'YarsRevengeDeterministic-v4',
704     'YarsRevengeNoFrameskip-v0',
705     'YarsRevengeNoFrameskip-v4',
706     'YarsRevenge-ram-v0',
707     'YarsRevenge-ram-v4',
708     'YarsRevenge-ramDeterministic-v0',
709     'YarsRevenge-ramDeterministic-v4',
710     'YarsRevenge-ramNoFrameskip-v0',
711     'YarsRevenge-ramNoFrameskip-v4',
712     'Zaxxon-v0',

```

(continues on next page)

(continued from previous page)

```

713         'Zaxxon-v4',
714         'ZaxxonDeterministic-v0',
715         'ZaxxonDeterministic-v4',
716         'ZaxxonNoFrameskip-v0',
717         'ZaxxonNoFrameskip-v4',
718         'Zaxxon-ram-v0',
719         'Zaxxon-ram-v4',
720         'Zaxxon-ramDeterministic-v0',
721         'Zaxxon-ramDeterministic-v4',
722         'Zaxxon-ramNoFrameskip-v0',
723         'Zaxxon-ramNoFrameskip-v4'],
724
725     # Classic control
726     'classic_control': [
727         'Acrobot-v1',
728         'CartPole-v1',
729         'CartPole-v0',
730         'MountainCar-v0',
731         'MountainCarContinuous-v0',
732         'Pendulum-v0'
733     ],
734
735     # Box2D
736     'box2d': [
737         'BipedalWalker-v2',
738         'BipedalWalkerHardcore-v2',
739         'CarRacing-v0',
740         'LunarLander-v2',
741         'LunarLanderContinuous-v2'
742     ],
743
744     # MuJoCo
745     'mujoco': [
746         'Ant-v2',
747         'HalfCheetah-v2',
748         'Hopper-v2',
749         'Humanoid-v2',
750         'HumanoidStandup-v2',
751         'InvertedDoublePendulum-v2',
752         'InvertedPendulum-v2',
753         'Reacher-v2',
754         'Swimmer-v2',
755         'Walker2d-v2'
756     ],
757
758     # Robotics
759     'robotics': [
760         'FetchPickAndPlace-v1',
761         'FetchPush-v1',
762         'FetchReach-v1',
763         'FetchSlide-v1',
764         'HandManipulateBlock-v0',
765         'HandManipulateEgg-v0',
766         'HandManipulatePen-v0',
767         'HandReach-v0'
768     ],
769

```

(continues on next page)

(continued from previous page)

```

770  ## Deepmind Control Suite  (need check!)
771  'dm_control': [
772      'AcrobotSparse-v0',
773      'BallinCupCatch-v0',
774      'CartpoleSwingup-v0',
775      'FingerTurn-v0',
776      'FishSwim-v0',
777      'CheetahRun-v0',
778      'HopperHop-v0',
779      'HumanoidStand-v0',
780      'HumanoidWalk-v0',
781      'HumanoidRun-v0',
782      'ManipulatorBringball-v0',
783      'PendulumSwingup-v0',
784      'Pointmass-v0',
785      'ReacherHard-v0',
786      'Swimmer-v0',
787      'WalkerRun-v0'
788  ],
789
790  ## RLBench
791  'rlbench': [
792      'BeatTheBuzz',
793      'BlockPyramid',
794      'ChangeChannel',
795      'ChangeClock',
796      'CloseBox',
797      'CloseDoor',
798      'CloseDrawer',
799      'CloseFridge',
800      'CloseGrill',
801      'CloseJar',
802      'CloseLaptopLid',
803      'CloseMicrowave',
804      'EmptyContainer',
805      'EmptyDishwasher',
806      'GetIceFromFridge',
807      'HangFrameOnHanger',
808      'HanoiSquare',
809      'HitBallWithQueue',
810      'Hockey',
811      'InsertUsbInComputer',
812      'LampOff',
813      'LampOn',
814      'LightBulbIn',
815      'LightBulbOut',
816      'MeatOffGrill',
817      'MeatOnGrill',
818      'MoveHanger',
819      'OpenBox',
820      'OpenDoor',
821      'OpenDrawer',
822      'OpenFridge',
823      'OpenGrill',
824      'OpenJar',
825      'OpenMicrowave',
826      'OpenOven',

```

(continues on next page)

(continued from previous page)

```

827     'OpenWindow',
828     'OpenWineBottle',
829     'PhoneOnBase',
830     'PickAndLift',
831     'PickUpCup',
832     'PlaceCups',
833     'PlaceHangerOnRack',
834     'PlaceShapeInShapeSorter',
835     'PlayJenga',
836     'PlugChargerInPowerSupply',
837     'PourFromCupToCup',
838     'PressSwitch',
839     'PushButton',
840     'PushButtons',
841     'PutBooksOnBookshelf',
842     'PutBottleInFridge',
843     'PutGroceriesInCupboard',
844     'PutItemInDrawer',
845     'PutKnifeInKnifeBlock',
846     'PutKnifeOnChoppingBoard',
847     'PutMoneyInSafe',
848     'PutPlateInColoredDishRack',
849     'PutRubbishInBin',
850     'PutShoesInBox',
851     'PutToiletRollOnStand',
852     'PutTrayInOven',
853     'PutUmbrellaInUmbrellaStand',
854     'ReachAndDrag',
855     'ReachTarget',
856     'RemoveCups',
857     'ScoopWithSpatula',
858     'ScrewNail',
859     'SetTheTable',
860     'SetupCheckers',
861     'SlideBlockToTarget',
862     'SlideCabinetOpen',
863     'SlideCabinetOpenAndPlaceCups',
864     'SolvePuzzle',
865     'StackBlocks',
866     'StackCups',
867     'StackWine',
868     'StraightenRope',
869     'SweepToDustpan',
870     'TakeCupOutFromCabinet',
871     'TakeFrameOffHanger',
872     'TakeItemOutOfDrawer',
873     'TakeLidOffSaucepan',
874     'TakeMoneyOutSafe',
875     'TakeOffWeighingScales',
876     'TakePlateOffColoredDishRack',
877     'TakeShoesOutOfBox',
878     'TakeToiletRollOffStand',
879     'TakeTrayOutOfOven',
880     'TakeUmbrellaOutOfUmbrellaStand',
881     'TakeUsbOutOfComputer',
882     'ToiletSeatDown',
883     'ToiletSeatUp',

```

(continues on next page)

(continued from previous page)

```
884     'TurnOvenOn',  
885     'TurnTap',  
886     'TvOff',  
887     'TvOn',  
888     'UnplugCharger',  
889     'WaterPlants',  
890     'WeighingScales',  
891     'WipeDesk'  
892 ]  
893 }
```


22.1 Math Utilities in RLzoo

Functions for mathematics utilization.

```
# Requirements tensorflow==2.0.0a0 tensorlayer==2.0.1
```

```
rlzoo.common.math_utils.flatten_dims(shapes)
```


23.1 Common Utilities in RLzoo

Functions for utilization.

Requirements tensorflow==2.0.0a0 tensorlayer==2.0.1

`rlzoo.common.utils.call_default_params(env, envtype, alg, default_seed=True)`

Get the default parameters for training from the default script

`rlzoo.common.utils.get_algorithm_module(algorithm, submodule)`

Get algorithm module in the corresponding folder

`rlzoo.common.utils.load_model(model, model_name, algorithm_name, env_name)`

load saved neural network model

Parameters

- **model** – tensorlayer.models.Model
- **model_name** – string, e.g. ‘model_sac_q1’
- **algorithm_name** – string, e.g. ‘SAC’

`rlzoo.common.utils.make_env(env_id)`

`rlzoo.common.utils.parse_all_args(parser)`

Parse known and unknown args

`rlzoo.common.utils.plot(episode_rewards, algorithm_name, env_name)`

plot the learning curve, saved as ./img/algorithm_name-env_name.png

Parameters

- **episode_rewards** – array of floats
- **algorithm_name** – string
- **env_name** – string

`rlzoo.common.utils.plot_save_log(episode_rewards, algorithm_name, env_name)`

plot the learning curve, saved as `./img/algorithm_name-env_name.png`, and save the rewards log as `./log/algorithm_name-env_name.npy`

Parameters

- **episode_rewards** – array of floats
- **algorithm_name** – string
- **env_name** – string

`rlzoo.common.utils.save_model(model, model_name, algorithm_name, env_name)`

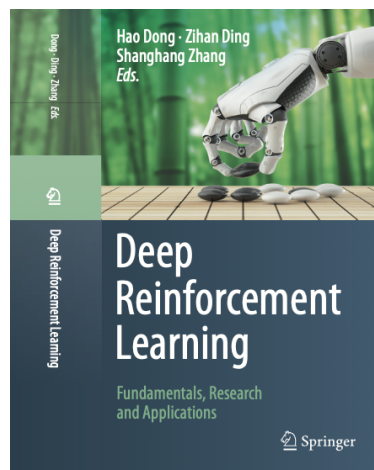
save trained neural network model

Parameters

- **model** – `tensorlayer.models.Model`
- **model_name** – string, e.g. `'model_sac_q1'`
- **algorithm_name** – string, e.g. `'SAC'`

`rlzoo.common.utils.set_seed(seed, env=None)`

set random seed for reproducibility



- You can get the [free PDF](#) if your institute has Springer license.

Deep reinforcement learning (DRL) relies on the intersection of reinforcement learning (RL) and deep learning (DL). It has been able to solve a wide range of complex decision-making tasks that were previously out of reach for a machine and famously contributed to the success of AlphaGo. Furthermore, it opens up numerous new applications in domains such as healthcare, robotics, smart grids, and finance.

Divided into three main parts, this book provides a comprehensive and self-contained introduction to DRL. The first part introduces the foundations of DL, RL and widely used DRL methods and discusses their implementation. The second part covers selected DRL research topics, which are useful for those wanting to specialize in DRL research. To help readers gain a deep understanding of DRL and quickly apply the techniques in practice, the third part presents mass applications, such as the intelligent transportation system and learning to run, with detailed explanations.

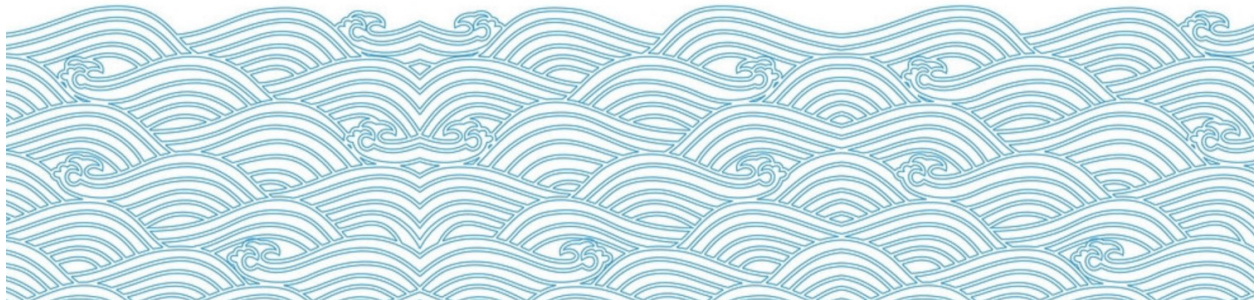
The book is intended for computer science students, both undergraduate and postgraduate, who would like to learn DRL from scratch, practice its implementation, and explore the research topics. This book also appeals to engineers and practitioners who do not have strong machine learning background, but want to quickly understand how DRL works and use the techniques in their applications.

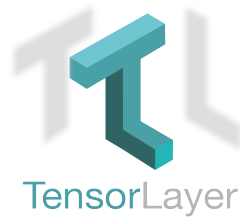
24.1 Editors

- Hao Dong - Peking University
- Zihan Ding - Princeton University
- Shanghang Zhang - University of California, Berkeley

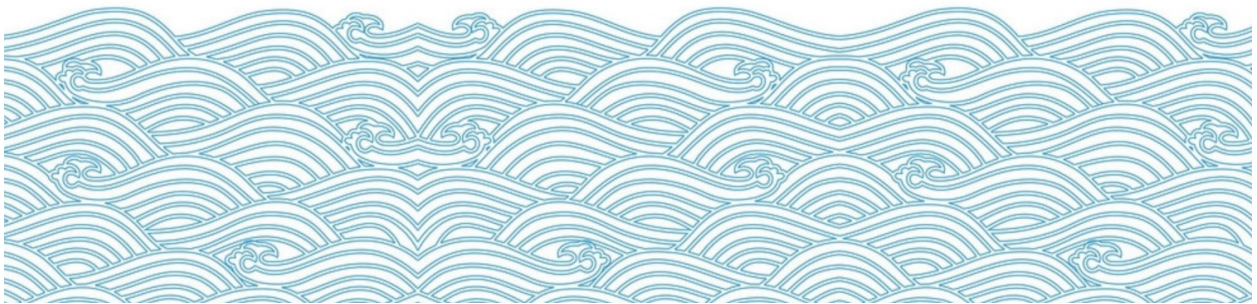
24.2 Authors

- Hao Dong - Peking University
- Zihan Ding - Princeton University
- Shanghang Zhang - University of California, Berkeley
- Hang Yuan - Oxford University
- Hongming Zhang - Peking University
- Jingqing Zhang - Imperial College London
- Yanhua Huang - Xiaohongshu Technology Co.
- Tianyang Yu - Nanchang University
- Huaqing Zhang - Google
- Ruitong Huang - Borealis AI





Different from RLzoo for simple usage with **high-level APIs**, the **RL tutorial** aims to make the reinforcement learning tutorial simple, transparent and straight-forward with **low-level APIs**, as this would not only benefits new learners of reinforcement learning, but also provide convenience for senior researchers to testify their new ideas quickly.



CHAPTER 26

Contributing

This project is under active development, if you want to join the core team, feel free to contact Zihan Ding at [zhd-ing\[at\]mail.ustc.edu.cn](mailto:zhd-ing[at]mail.ustc.edu.cn)

CHAPTER 27

Citation

- genindex
- modindex
- search



r

- `rlzoo.algorithms.a3c.default`, 22
- `rlzoo.algorithms.ac.default`, 20
- `rlzoo.algorithms.ddpg.default`, 25
- `rlzoo.algorithms.dppo.default`, 49
- `rlzoo.algorithms.dqn.default`, 14
- `rlzoo.algorithms.pg.default`, 17
- `rlzoo.algorithms.ppo.default`, 43
- `rlzoo.algorithms.sac.default`, 33
- `rlzoo.algorithms.td3.default`, 29
- `rlzoo.algorithms.trpo.default`, 38
- `rlzoo.common.basic_nets`, 51
- `rlzoo.common.buffer`, 61
- `rlzoo.common.distributions`, 65
- `rlzoo.common.env_list`, 71
- `rlzoo.common.env_wrappers`, 67
- `rlzoo.common.math_utils`, 89
- `rlzoo.common.utils`, 91

Symbols

`__call__()` (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` method), 55
`__call__()` (`rlzoo.common.policy_networks.StochasticPolicyNetwork` method), 56
`__call__()` (`rlzoo.common.value_networks.QNetwork` method), 58
`__call__()` (`rlzoo.common.value_networks.ValueNetwork` method), 57
`__dict__` (`rlzoo.common.buffer.ReplayBuffer` attribute), 62
`__dict__` (`rlzoo.common.buffer.SegmentTree` attribute), 63
`__getitem__()` (`rlzoo.common.buffer.SegmentTree` method), 63
`__init__()` (`rlzoo.common.buffer.HindsightReplayBuffer` method), 61
`__init__()` (`rlzoo.common.buffer.MinSegmentTree` method), 62
`__init__()` (`rlzoo.common.buffer.PrioritizedReplayBuffer` method), 62
`__init__()` (`rlzoo.common.buffer.ReplayBuffer` method), 62
`__init__()` (`rlzoo.common.buffer.SegmentTree` method), 63
`__init__()` (`rlzoo.common.buffer.SumSegmentTree` method), 63
`__init__()` (`rlzoo.common.policy_networks.DeterministicContinuousPolicyNetwork` method), 54
`__init__()` (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` method), 55
`__init__()` (`rlzoo.common.policy_networks.StochasticContinuousPolicyNetwork` method), 53
`__init__()` (`rlzoo.common.policy_networks.StochasticPolicyNetwork` method), 56
`__init__()` (`rlzoo.common.value_networks.MlpQNetwork` method), 58
`__init__()` (`rlzoo.common.value_networks.QNetwork` method), 58
`__init__()` (`rlzoo.common.value_networks.ValueNetwork` method), 57
`__len__()` (`rlzoo.common.buffer.ReplayBuffer` method), 63
`__module__` (`rlzoo.common.buffer.HindsightReplayBuffer` attribute), 61
`__module__` (`rlzoo.common.buffer.MinSegmentTree` attribute), 62
`__module__` (`rlzoo.common.buffer.PrioritizedReplayBuffer` attribute), 62
`__module__` (`rlzoo.common.buffer.ReplayBuffer` attribute), 63
`__module__` (`rlzoo.common.buffer.SegmentTree` attribute), 63
`__module__` (`rlzoo.common.buffer.SumSegmentTree` attribute), 64
`__module__` (`rlzoo.common.policy_networks.DeterministicContinuousPolicyNetwork` attribute), 54
`__module__` (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` attribute), 55
`__module__` (`rlzoo.common.policy_networks.StochasticContinuousPolicyNetwork` attribute), 54
`__module__` (`rlzoo.common.policy_networks.StochasticPolicyNetwork` attribute), 56
`__module__` (`rlzoo.common.value_networks.MlpQNetwork` attribute), 58
`__module__` (`rlzoo.common.value_networks.QNetwork` attribute), 59
`__module__` (`rlzoo.common.value_networks.ValueNetwork` attribute), 57
`__setitem__()` (`rlzoo.common.buffer.SegmentTree` method), 63
`__weakref__` (`rlzoo.common.buffer.ReplayBuffer` attribute), 63
`__weakref__` (`rlzoo.common.buffer.SegmentTree` attribute), 63

A

A3C (class in `rlzoo.algorithms.a3c.a3c`), 22

`a_train()` (`rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP` method), 48
`a_train()` (`rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY` method), 46
`a_train()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 41
`a_train()` (`rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY` method), 40
`a_train()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 36
`AC` (class in `rlzoo.algorithms.ac.ac`), 20
`action_shape` (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` attribute), 55
`action_shape` (`rlzoo.common.policy_networks.StochasticPolicyNetwork` attribute), 56
`action_shape` (`rlzoo.common.value_networks.QNetwork` attribute), 59
`action_space` (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` attribute), 55
`action_space` (`rlzoo.common.policy_networks.StochasticPolicyNetwork` attribute), 56
`action_space` (`rlzoo.common.value_networks.QNetwork` attribute), 59
`assign_params_from_flat()` (`rlzoo.algorithms.trpo.trpo.TRPO` static method), 36
`atari()` (in module `rlzoo.algorithms.a3c.default`), 22
`atari()` (in module `rlzoo.algorithms.ac.default`), 20
`atari()` (in module `rlzoo.algorithms.dppo.default`), 49
`atari()` (in module `rlzoo.algorithms.dqn.default`), 14
`atari()` (in module `rlzoo.algorithms.pg.default`), 17
`atari()` (in module `rlzoo.algorithms.ppo.default`), 43
`atari()` (in module `rlzoo.algorithms.trpo.default`), 38
B
`box2d()` (in module `rlzoo.algorithms.a3c.default`), 22
`box2d()` (in module `rlzoo.algorithms.ac.default`), 20
`box2d()` (in module `rlzoo.algorithms.ddpg.default`), 25
`box2d()` (in module `rlzoo.algorithms.dppo.default`), 49
`box2d()` (in module `rlzoo.algorithms.pg.default`), 17
`box2d()` (in module `rlzoo.algorithms.ppo.default`), 43
`box2d()` (in module `rlzoo.algorithms.sac.default`), 33
`box2d()` (in module `rlzoo.algorithms.td3.default`), 29
`box2d()` (in module `rlzoo.algorithms.trpo.default`), 38
`build_env()` (in module `rlzoo.common.env_wrappers`), 67
C
`c_train()` (`rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP` method), 48
`c_train()` (`rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY` method), 46
`c_train()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 42
D
`DDPG` (class in `rlzoo.algorithms.ddpg.ddpg`), 24
`ContinuousPolicyNetwork` (class in `rlzoo.common.policy_networks`), 54
`DeterministicPolicyNetwork` (class in `rlzoo.common.policy_networks`), 54
`cal_adv()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 42
`cal_adv()` (`rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY` method), 46
`cal_adv()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 36
`call_default_params()` (in module `rlzoo.common.utils`), 91
`Categorical` (class in `rlzoo.common.distributions`), 6
`cg()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 36
`control()` (in module `rlzoo.algorithms.a3c.default`), 22
`classic_control()` (in module `rlzoo.algorithms.ac.default`), 20
`classic_control()` (in module `rlzoo.algorithms.ddpg.default`), 25
`classic_control()` (in module `rlzoo.algorithms.dppo.default`), 49
`classic_control()` (in module `rlzoo.algorithms.dqn.default`), 14
`classic_control()` (in module `rlzoo.algorithms.pg.default`), 17
`classic_control()` (in module `rlzoo.algorithms.ppo.default`), 43
`classic_control()` (in module `rlzoo.algorithms.sac.default`), 33
`classic_control()` (in module `rlzoo.algorithms.td3.default`), 29
`classic_control()` (in module `rlzoo.algorithms.trpo.default`), 38
`ClipRewardEnv` (class in `rlzoo.common.env_wrappers`), 69
`close()` (`rlzoo.common.env_wrappers.SubprocVecEnv` method), 69
`CNN()` (in module `rlzoo.common.basic_nets`), 51
`CNNModel()` (in module `rlzoo.common.basic_nets`), 51
`CreateInputLayer()` (in module `rlzoo.common.basic_nets`), 51

DiagGaussian (class in *rlzoo.common.distributions*), 65
 Distribution (class in *rlzoo.common.distributions*), 66
 dm_control() (in module *rl-zoo.algorithms.a3c.default*), 22
 dm_control() (in module *rl-zoo.algorithms.ac.default*), 20
 dm_control() (in module *rl-zoo.algorithms.ddpg.default*), 25
 dm_control() (in module *rl-zoo.algorithms.dppo.default*), 50
 dm_control() (in module *rl-zoo.algorithms.pg.default*), 17
 dm_control() (in module *rl-zoo.algorithms.ppo.default*), 43
 dm_control() (in module *rl-zoo.algorithms.sac.default*), 33
 dm_control() (in module *rl-zoo.algorithms.td3.default*), 29
 dm_control() (in module *rl-zoo.algorithms.trpo.default*), 38
 DmObsTrans (class in *rlzoo.common.env_wrappers*), 70
 DPPO_CLIP (class in *rl-zoo.algorithms.dppo_clip.dppo_clip*), 48
 DPPO_PENALTY (class in *rl-zoo.algorithms.dppo_penalty.dppo_penalty*), 46
 DQN (class in *rlzoo.algorithms.dqn.dqn*), 13
E
 ema_update() (*rlzoo.algorithms.ddpg.ddpg.DDPG* method), 24
 entropy() (*rlzoo.common.distributions.Categorical* method), 65
 entropy() (*rlzoo.common.distributions.DiagGaussian* method), 66
 entropy() (*rlzoo.common.distributions.Distribution* method), 66
 EpisodicLifeEnv (class in *rl-zoo.common.env_wrappers*), 68
 eval() (*rlzoo.algorithms.trpo.trpo.TRPO* method), 36
 evaluate() (*rlzoo.algorithms.sac.sac.SAC* method), 32
 evaluate() (*rlzoo.algorithms.td3.td3.TD3* method), 28
 expand_dims() (in module *rl-zoo.common.distributions*), 66
F
 find_prefixsum_idx() (*rl-zoo.common.buffer.SumSegmentTree* method), 64
 FireResetEnv (class in *rl-zoo.common.env_wrappers*), 68
 flat_concat() (*rlzoo.algorithms.trpo.trpo.TRPO* static method), 36
 flatten_dims() (in module *rl-zoo.common.math_utils*), 89
 FrameStack (class in *rlzoo.common.env_wrappers*), 69
G
 get_action() (*rlzoo.algorithms.ac.ac.AC* method), 20
 get_action() (*rlzoo.algorithms.ddpg.ddpg.DDPG* method), 24
 get_action() (*rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP* method), 48
 get_action() (*rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY* method), 46
 get_action() (*rlzoo.algorithms.dqn.dqn.DQN* method), 13
 get_action() (*rlzoo.algorithms.pg.pg.PG* method), 16
 get_action() (*rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP* method), 42
 get_action() (*rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY* method), 40
 get_action() (*rlzoo.algorithms.sac.sac.SAC* method), 32
 get_action() (*rlzoo.algorithms.td3.td3.TD3* method), 28
 get_action() (*rlzoo.algorithms.trpo.trpo.TRPO* method), 36
 get_action_greedy() (*rlzoo.algorithms.ac.ac.AC* method), 20
 get_action_greedy() (*rl-zoo.algorithms.ddpg.ddpg.DDPG* method), 24
 get_action_greedy() (*rl-zoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP* method), 48
 get_action_greedy() (*rl-zoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY* method), 46
 get_action_greedy() (*rl-zoo.algorithms.dqn.dqn.DQN* method), 13
 get_action_greedy() (*rlzoo.algorithms.pg.pg.PG* method), 16
 get_action_greedy() (*rl-zoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP* method), 42
 get_action_greedy() (*rl-zoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY* method), 40

`get_action_greedy()` (`rlzoo.algorithms.sac.sac.SAC` method), 32
`get_action_greedy()` (`rlzoo.algorithms.td3.td3.TD3` method), 28
`get_action_greedy()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 36
`get_algorithm_module()` (in module `rlzoo.common.utils`), 91
`get_envlist()` (in module `rlzoo.common.env_list`), 71
`get_param()` (`rlzoo.common.distributions.Categorical` method), 65
`get_param()` (`rlzoo.common.distributions.DiagGaussian` method), 66
`get_pi_params()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 36
`get_v()` (`rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP` method), 48
`get_v()` (`rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY` method), 46
`get_v()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 42
`get_v()` (`rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY` method), 40
`get_v()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
`GOAL_EPISODE` (`rlzoo.common.buffer.HindsightReplayBuffer` attribute), 61
`GOAL_FUTURE` (`rlzoo.common.buffer.HindsightReplayBuffer` attribute), 61
`GOAL_RANDOM` (`rlzoo.common.buffer.HindsightReplayBuffer` attribute), 61
`gradient()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
`greedy_sample()` (`rlzoo.common.distributions.Categorical` method), 65
`greedy_sample()` (`rlzoo.common.distributions.DiagGaussian` method), 66

H

`hessian_vector_product()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
`HindsightReplayBuffer` (class in `rlzoo.common.buffer`), 61

K

`kl()` (`rlzoo.common.distributions.Categorical` method), 65
`kl()` (`rlzoo.common.distributions.DiagGaussian` method), 66

L

`kl()` (`rlzoo.common.distributions.Distribution` method), 66
`LazyFrames` (class in `rlzoo.common.env_wrappers`), 69
`learn()` (`rlzoo.algorithms.a3c.a3c.A3C` method), 22
`learn()` (`rlzoo.algorithms.ac.ac.AC` method), 20
`learn()` (`rlzoo.algorithms.ddpg.ddpg.DDPG` method), 24
`learn()` (`rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP` method), 48
`learn()` (`rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY` method), 47
`learn()` (`rlzoo.algorithms.dqn.dqn.DQN` method), 14
`learn()` (`rlzoo.algorithms.pg.pg.PG` method), 16
`learn()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 42
`learn()` (`rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY` method), 40
`learn()` (`rlzoo.algorithms.sac.sac.SAC` method), 32
`learn()` (`rlzoo.algorithms.td3.td3.TD3` method), 28
`learn()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
`load_ckpt()` (`rlzoo.algorithms.ac.ac.AC` method), 20
`load_ckpt()` (`rlzoo.algorithms.ddpg.ddpg.DDPG` method), 24
`load_ckpt()` (`rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP` method), 49
`load_ckpt()` (`rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY` method), 47
`load_ckpt()` (`rlzoo.algorithms.dqn.dqn.DQN` method), 14
`load_ckpt()` (`rlzoo.algorithms.pg.pg.PG` method), 16
`load_ckpt()` (`rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP` method), 43
`load_ckpt()` (`rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY` method), 41
`load_ckpt()` (`rlzoo.algorithms.sac.sac.SAC` method), 32
`load_ckpt()` (`rlzoo.algorithms.td3.td3.TD3` method), 28
`load_ckpt()` (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
`load_model()` (in module `rlzoo.common.utils`), 91
`logp()` (`rlzoo.common.distributions.Categorical` method), 65
`logp()` (`rlzoo.common.distributions.DiagGaussian` method), 66
`logp()` (`rlzoo.common.distributions.Distribution` method), 66

M

`make_dist()` (in module `rlzoo.common.distributions`), 66

[make_env\(\)](#) (in module `rlzoo.common.utils`), 91
[MaxAndSkipEnv](#) (class in `rlzoo.common.env_wrappers`), 69
[min\(\)](#) (`rlzoo.common.buffer.MinSegmentTree` method), 62
[MinSegmentTree](#) (class in `rlzoo.common.buffer`), 61
[MLP\(\)](#) (in module `rlzoo.common.basic_nets`), 51
[MLPModel\(\)](#) (in module `rlzoo.common.basic_nets`), 52
[MlpQNetwork](#) (class in `rlzoo.common.value_networks`), 58
[Monitor](#) (class in `rlzoo.common.env_wrappers`), 70
[mujoco\(\)](#) (in module `rlzoo.algorithms.a3c.default`), 22
[mujoco\(\)](#) (in module `rlzoo.algorithms.ac.default`), 20
[mujoco\(\)](#) (in module `rlzoo.algorithms.ddpg.default`), 25
[mujoco\(\)](#) (in module `rlzoo.algorithms.dppo.default`), 50
[mujoco\(\)](#) (in module `rlzoo.algorithms.pg.default`), 17
[mujoco\(\)](#) (in module `rlzoo.algorithms.ppo.default`), 43
[mujoco\(\)](#) (in module `rlzoo.algorithms.sac.default`), 33
[mujoco\(\)](#) (in module `rlzoo.algorithms.td3.default`), 29
[mujoco\(\)](#) (in module `rlzoo.algorithms.trpo.default`), 38

N

[ndim](#) (`rlzoo.common.distributions.Categorical` attribute), 65
[ndim](#) (`rlzoo.common.distributions.DiagGaussian` attribute), 66
[neglogp\(\)](#) (`rlzoo.common.distributions.Categorical` method), 65
[neglogp\(\)](#) (`rlzoo.common.distributions.DiagGaussian` method), 66
[neglogp\(\)](#) (`rlzoo.common.distributions.Distribution` method), 66
[NoopResetEnv](#) (class in `rlzoo.common.env_wrappers`), 68
[NormalizedActions](#) (class in `rlzoo.common.env_wrappers`), 70

O

[observation\(\)](#) (`rlzoo.common.env_wrappers.WarpFrame` method), 69

P

[parse_all_args\(\)](#) (in module `rlzoo.common.utils`), 91
[PG](#) (class in `rlzoo.algorithms.pg.pg`), 16
[pi_loss\(\)](#) (`rlzoo.algorithms.trpo.trpo.TRPO` method), 37
[plot\(\)](#) (in module `rlzoo.common.utils`), 91
[plot_save_log\(\)](#) (in module `rlzoo.common.utils`), 91

[PPO_CLIP](#) (class in `rlzoo.algorithms.ppo_clip.ppo_clip`), 41
[PPO_PENALTY](#) (class in `rlzoo.algorithms.ppo_penalty.ppo_penalty`), 40
[PrioritizedReplayBuffer](#) (class in `rlzoo.common.buffer`), 62
[push\(\)](#) (`rlzoo.common.buffer.HindsightReplayBuffer` method), 61
[push\(\)](#) (`rlzoo.common.buffer.PrioritizedReplayBuffer` method), 62
[push\(\)](#) (`rlzoo.common.buffer.ReplayBuffer` method), 63
[push_episode\(\)](#) (`rlzoo.common.buffer.HindsightReplayBuffer` method), 61

Q

[QNetwork](#) (class in `rlzoo.common.value_networks`), 58

R

[random_sample\(\)](#) (`rlzoo.common.policy_networks.DeterministicPolicyNetwork` method), 55
[random_sample\(\)](#) (`rlzoo.common.policy_networks.StochasticPolicyNetwork` method), 56
[reduce\(\)](#) (`rlzoo.common.buffer.SegmentTree` method), 63
[ReplayBuffer](#) (class in `rlzoo.common.buffer`), 62
[reset\(\)](#) (`rlzoo.common.env_wrappers.DmObsTrans` method), 70
[reset\(\)](#) (`rlzoo.common.env_wrappers.EpisodicLifeEnv` method), 68
[reset\(\)](#) (`rlzoo.common.env_wrappers.FireResetEnv` method), 68
[reset\(\)](#) (`rlzoo.common.env_wrappers.FrameStack` method), 69
[reset\(\)](#) (`rlzoo.common.env_wrappers.MaxAndSkipEnv` method), 69
[reset\(\)](#) (`rlzoo.common.env_wrappers.Monitor` method), 70
[reset\(\)](#) (`rlzoo.common.env_wrappers.NoopResetEnv` method), 68
[reset\(\)](#) (`rlzoo.common.env_wrappers.SubprocVecEnv` method), 70
[reset\(\)](#) (`rlzoo.common.env_wrappers.TimeLimit` method), 67
[reset\(\)](#) (`rlzoo.common.env_wrappers.VecFrameStack` method), 70
[reward\(\)](#) (`rlzoo.common.env_wrappers.ClipRewardEnv` method), 69
[reward\(\)](#) (`rlzoo.common.env_wrappers.RewardShaping` method), 69

RewardShaping (class in rlzoo.common.env_wrappers), 69
 rlbench() (in module rlzoo.algorithms.a3c.default), 22
 rlbench() (in module rlzoo.algorithms.ac.default), 20
 rlbench() (in module rlzoo.algorithms.ddpg.default), 25
 rlbench() (in module rlzoo.algorithms.pg.default), 17
 rlbench() (in module rlzoo.algorithms.ppo.default), 43
 rlbench() (in module rlzoo.algorithms.sac.default), 33
 rlbench() (in module rlzoo.algorithms.td3.default), 29
 rlbench() (in module rlzoo.algorithms.trpo.default), 38
 rlzoo.algorithms.a3c.default (module), 22
 rlzoo.algorithms.ac.default (module), 20
 rlzoo.algorithms.ddpg.default (module), 25
 rlzoo.algorithms.dppo.default (module), 49
 rlzoo.algorithms.dqn.default (module), 14
 rlzoo.algorithms.pg.default (module), 17
 rlzoo.algorithms.ppo.default (module), 43
 rlzoo.algorithms.sac.default (module), 33
 rlzoo.algorithms.td3.default (module), 29
 rlzoo.algorithms.trpo.default (module), 38
 rlzoo.common.basic_nets (module), 51
 rlzoo.common.buffer (module), 61
 rlzoo.common.distributions (module), 65
 rlzoo.common.env_list (module), 71
 rlzoo.common.env_wrappers (module), 67
 rlzoo.common.math_utils (module), 89
 rlzoo.common.utils (module), 91
 robotics() (in module rlzoo.algorithms.a3c.default), 22
 robotics() (in module rlzoo.algorithms.ac.default), 20
 robotics() (in module rlzoo.algorithms.ddpg.default), 25
 robotics() (in module rlzoo.algorithms.dppo.default), 50
 robotics() (in module rlzoo.algorithms.pg.default), 17
 robotics() (in module rlzoo.algorithms.ppo.default), 43
 robotics() (in module rlzoo.algorithms.sac.default), 33
 robotics() (in module rlzoo.algorithms.td3.default), 29
 robotics() (in module rlzoo.algorithms.trpo.default), 38
 S
 SAC (class in rlzoo.algorithms.sac.sac), 32
 sample() (rlzoo.common.buffer.PrioritizedReplayBuffer method), 62
 sample() (rlzoo.common.buffer.ReplayBuffer method), 63
 sample() (rlzoo.common.distributions.Categorical method), 65
 sample() (rlzoo.common.distributions.DiagGaussian method), 66
 sample() (rlzoo.common.distributions.Distribution method), 66
 sample_action() (rlzoo.algorithms.ddpg.ddpg.DDPG method), 24
 sample_action() (rlzoo.algorithms.sac.sac.SAC method), 32
 sample_action() (rlzoo.algorithms.td3.td3.TD3 method), 28
 save_ckpt() (rlzoo.algorithms.ac.ac.AC method), 20
 save_ckpt() (rlzoo.algorithms.ddpg.ddpg.DDPG method), 24
 save_ckpt() (rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP method), 49
 save_ckpt() (rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALTY method), 47
 save_ckpt() (rlzoo.algorithms.dqn.dqn.DQN method), 14
 save_ckpt() (rlzoo.algorithms.pg.pg.PG method), 16
 save_ckpt() (rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP method), 43
 save_ckpt() (rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY method), 41
 save_ckpt() (rlzoo.algorithms.sac.sac.SAC method), 32
 save_ckpt() (rlzoo.algorithms.td3.td3.TD3 method), 28
 save_ckpt() (rlzoo.algorithms.trpo.trpo.TRPO method), 37
 save_model() (in module rlzoo.common.utils), 92
 SegmentTree (class in rlzoo.common.buffer), 63
 set_param() (rlzoo.common.distributions.Categorical method), 65
 set_param() (rlzoo.common.distributions.DiagGaussian method), 66
 set_param() (rlzoo.common.distributions.Distribution method), 66
 set_pi_params() (rlzoo.algorithms.trpo.trpo.TRPO method), 38
 set_seed() (in module rlzoo.common.utils), 92
 state_shape (rlzoo.common.policy_networks.DeterministicPolicyNetwork attribute), 55
 state_shape (rlzoo.common.policy_networks.StochasticPolicyNetwork attribute), 56
 state_shape (rlzoo.common.value_networks.QNetwork attribute), 59

`state_shape (rlzoo.common.value_networks.ValueNetwork`
`attribute), 57`
`state_space (rlzoo.common.policy_networks.DeterministicPolicyNetwork`
`attribute), 55`
`state_space (rlzoo.common.policy_networks.StochasticPolicyNetwork`
`attribute), 56`
`state_space (rlzoo.common.value_networks.QNetwork`
`attribute), 59`
`state_space (rlzoo.common.value_networks.ValueNetwork`
`attribute), 58`
`step () (rlzoo.common.env_wrappers.DmObsTrans`
`method), 70`
`step () (rlzoo.common.env_wrappers.EpisodicLifeEnv`
`method), 68`
`step () (rlzoo.common.env_wrappers.FireResetEnv`
`method), 68`
`step () (rlzoo.common.env_wrappers.FrameStack`
`method), 69`
`step () (rlzoo.common.env_wrappers.MaxAndSkipEnv`
`method), 69`
`step () (rlzoo.common.env_wrappers.Monitor method),`
`70`
`step () (rlzoo.common.env_wrappers.NoopResetEnv`
`method), 68`
`step () (rlzoo.common.env_wrappers.SubprocVecEnv`
`method), 70`
`step () (rlzoo.common.env_wrappers.TimeLimit`
`method), 67`
`step () (rlzoo.common.env_wrappers.VecFrameStack`
`method), 70`
`StochasticContinuousPolicyNetwork (class`
`in rlzoo.common.policy_networks), 53`
`StochasticPolicyNetwork (class in rl-`
`zoo.common.policy_networks), 55`
`store_transition () (rl-`
`zoo.algorithms.ddpg.ddpg.DDPG method),`
`25`
`store_transition () (rl-`
`zoo.algorithms.dqn.dqn.DQN method), 14`
`store_transition () (rlzoo.algorithms.pg.pg.PG`
`method), 16`
`SubprocVecEnv (class in rl-`
`zoo.common.env_wrappers), 69`
`sum () (rlzoo.common.buffer.SumSegmentTree method),`
`64`
`SumSegmentTree (class in rlzoo.common.buffer), 63`
`sync () (rlzoo.algorithms.dqn.dqn.DQN method), 14`

T

`target_ini () (rlzoo.algorithms.sac.sac.SAC`
`method), 32`
`target_ini () (rlzoo.algorithms.td3.td3.TD3`
`method), 28`

`target_soft_update () (rl-`
`zoo.algorithms.sac.sac.SAC method), 32`
`update () (rlzoo.algorithms.td3.td3.TD3 method), 28`
`PolicyNetwork (class in rlzoo.common.env_wrappers), 67`
`TimeLimit (class in rlzoo.common.env_wrappers), 67`
`TRPO (class in rlzoo.algorithms.trpo.trpo), 36`

U

`update () (rlzoo.algorithms.ac.ac.AC method), 20`
`update () (rlzoo.algorithms.ddpg.ddpg.DDPG`
`method), 25`
`update () (rlzoo.algorithms.dppo_clip.dppo_clip.DPPO_CLIP`
`method), 49`
`update () (rlzoo.algorithms.dppo_penalty.dppo_penalty.DPPO_PENALT`
`method), 47`
`update () (rlzoo.algorithms.dqn.dqn.DQN method), 14`
`update () (rlzoo.algorithms.pg.pg.PG method), 16`
`update () (rlzoo.algorithms.ppo_clip.ppo_clip.PPO_CLIP`
`method), 43`
`update () (rlzoo.algorithms.ppo_penalty.ppo_penalty.PPO_PENALTY`
`method), 41`
`update () (rlzoo.algorithms.sac.sac.SAC method), 32`
`update () (rlzoo.algorithms.td3.td3.TD3 method), 28`
`update () (rlzoo.algorithms.trpo.trpo.TRPO method),`
`38`
`update_priorities () (rl-`
`zoo.common.buffer.PrioritizedReplayBuffer`
`method), 62`

V

`ValueNetwork (class in rl-`
`zoo.common.value_networks), 57`
`VecFrameStack (class in rl-`
`zoo.common.env_wrappers), 70`

W

`WarpFrame (class in rlzoo.common.env_wrappers), 69`