
Tianshou

Release 0.2.3

Tianshou contributors

Jun 01, 2020

TUTORIALS

1	Installation	3
2	Indices and tables	43
	Bibliography	45
	Python Module Index	47
	Index	49

Tianshou () is a reinforcement learning platform based on pure PyTorch. Unlike existing reinforcement learning libraries, which are mainly based on TensorFlow, have many nested classes, unfriendly API, or slow-speed, Tianshou provides a fast-speed framework and pythonic API for building the deep reinforcement learning agent. The supported interface algorithms include:

- *PGPolicy* Policy Gradient
- *DQNPolicy* Deep Q-Network
- *DQNPolicy* Double DQN with n-step returns
- *A2CPolicy* Advantage Actor-Critic
- *DDPGPolicy* Deep Deterministic Policy Gradient
- *PPOPolicy* Proximal Policy Optimization
- *TD3Policy* Twin Delayed DDPG
- *SACPolicy* Soft Actor-Critic
- *ImitationPolicy* Imitation Learning
- *PrioritizedReplayBuffer* Prioritized Experience Replay
- *compute_episodic_return()* Generalized Advantage Estimator

Tianshou supports parallel workers for all algorithms as well. All of these algorithms are reformatted as replay-buffer based algorithms.

INSTALLATION

Tianshou is currently hosted on [PyPI](#). You can simply install Tianshou with the following command (with Python \geq 3.6):

```
pip3 install tianshou
```

You can also install with the newest version through GitHub:

```
pip3 install git+https://github.com/thu-ml/tianshou.git@master
```

If you use Anaconda or Miniconda, you can install Tianshou through the following command lines:

```
# create a new virtualenv and install pip, change the env name if you like
conda create -n myenv pip
# activate the environment
conda activate myenv
# install tianshou
pip install tianshou
```

After installation, open your python console and type

```
import tianshou as ts
print(ts.__version__)
```

If no error occurs, you have successfully installed Tianshou.

Tianshou is still under development, you can also check out the documents in stable version through tianshou.readthedocs.io/en/stable/.

1.1 Deep Q Network

Deep reinforcement learning has achieved significant successes in various applications. **Deep Q Network** (DQN) [MKS+15] is the pioneer one. In this tutorial, we will show how to train a DQN agent on CartPole with Tianshou step by step. The full script is at [test/discrete/test_dqn.py](#).

Contrary to existing Deep RL libraries such as [RLlib](#), which could only accept a config specification of hyperparameters, network, and others, Tianshou provides an easy way of construction through the code-level.

1.1.1 Make an Environment

First of all, you have to make an environment for your agent to interact with. For environment interfaces, we follow the convention of [OpenAI Gym](#). In your Python code, simply import Tianshou and make the environment:

```
import gym
import tianshou as ts

env = gym.make('CartPole-v0')
```

CartPole-v0 is a simple environment with a discrete action space, for which DQN applies. You have to identify whether the action space is continuous or discrete and apply eligible algorithms. DDPG [LHP+16], for example, could only be applied to continuous action spaces, while almost all other policy gradient methods could be applied to both, depending on the probability distribution on the action.

1.1.2 Setup Multi-environment Wrapper

It is available if you want the original `gym.Env`:

```
train_envs = gym.make('CartPole-v0')
test_envs = gym.make('CartPole-v0')
```

Tianshou supports parallel sampling for all algorithms. It provides three types of vectorized environment wrapper: *VectorEnv*, *SubprocVectorEnv*, and *RayVectorEnv*. It can be used as follows:

```
train_envs = ts.env.VectorEnv([lambda: gym.make('CartPole-v0') for _ in range(8)])
test_envs = ts.env.VectorEnv([lambda: gym.make('CartPole-v0') for _ in range(100)])
```

Here, we set up 8 environments in `train_envs` and 100 environments in `test_envs`.

For the demonstration, here we use the second block of codes.

1.1.3 Build the Network

Tianshou supports any user-defined PyTorch networks and optimizers but with the limitation of input and output API. Here is an example code:

```
import torch, numpy as np
from torch import nn

class Net(nn.Module):
    def __init__(self, state_shape, action_shape):
        super().__init__()
        self.model = nn.Sequential(*[
            nn.Linear(np.prod(state_shape), 128), nn.ReLU(inplace=True),
            nn.Linear(128, 128), nn.ReLU(inplace=True),
            nn.Linear(128, 128), nn.ReLU(inplace=True),
            nn.Linear(128, np.prod(action_shape))
        ])
    def forward(self, obs, state=None, info={}):
        if not isinstance(obs, torch.Tensor):
            obs = torch.tensor(obs, dtype=torch.float)
        batch = obs.shape[0]
        logits = self.model(obs.view(batch, -1))
        return logits, state
```

(continues on next page)

(continued from previous page)

```

state_shape = env.observation_space.shape or env.observation_space.n
action_shape = env.action_space.shape or env.action_space.n
net = Net(state_shape, action_shape)
optim = torch.optim.Adam(net.parameters(), lr=1e-3)

```

The rules of self-defined networks are:

1. Input: observation `obs` (may be a `numpy.ndarray` or `torch.Tensor`), hidden state `state` (for RNN usage), and other information `info` provided by the environment.
2. Output: some logits and the next hidden state `state`. The logits could be a tuple instead of a `torch.Tensor`. It depends on how the policy process the network output. For example, in PPO [SWD+17], the return of the network might be `(mu, sigma), state` for Gaussian policy.

1.1.4 Setup Policy

We use the defined `net` and `optim`, with extra policy hyper-parameters, to define a policy. Here we define a DQN policy with using a target network:

```

policy = ts.policy.DQNPolicy(net, optim,
    discount_factor=0.9, estimation_step=3,
    use_target_network=True, target_update_freq=320)

```

1.1.5 Setup Collector

The collector is a key concept in Tianshou. It allows the policy to interact with different types of environments conveniently. In each step, the collector will let the policy perform (at least) a specified number of steps or episodes and store the data in a replay buffer.

```

train_collector = ts.data.Collector(policy, train_envs, ts.data.
    ↳ReplayBuffer(size=20000))
test_collector = ts.data.Collector(policy, test_envs)

```

1.1.6 Train Policy with a Trainer

Tianshou provides `onpolicy_trainer` and `offpolicy_trainer`. The trainer will automatically stop training when the policy reach the stop condition `stop_fn` on test collector. Since DQN is an off-policy algorithm, we use the `offpolicy_trainer` as follows:

```

result = ts.trainer.offpolicy_trainer(
    policy, train_collector, test_collector,
    max_epoch=10, step_per_epoch=1000, collect_per_step=10,
    episode_per_test=100, batch_size=64,
    train_fn=lambda e: policy.set_eps(0.1),
    test_fn=lambda e: policy.set_eps(0.05),
    stop_fn=lambda x: x >= env.spec.reward_threshold,
    writer=None)
print(f'Finished training! Use {result["duration"]}')

```

The meaning of each parameter is as follows:

- `max_epoch`: The maximum of epochs for training. The training process might be finished before reaching the `max_epoch`;
- `step_per_epoch`: The number of step for updating policy network in one epoch;
- `collect_per_step`: The number of frames the collector would collect before the network update. For example, the code above means “collect 10 frames and do one policy network update”;
- `episode_per_test`: The number of episodes for one policy evaluation.
- `batch_size`: The batch size of sample data, which is going to feed in the policy network.
- `train_fn`: A function receives the current number of epoch index and performs some operations at the beginning of training in this epoch. For example, the code above means “reset the epsilon to 0.1 in DQN before training”.
- `test_fn`: A function receives the current number of epoch index and performs some operations at the beginning of testing in this epoch. For example, the code above means “reset the epsilon to 0.05 in DQN before testing”.
- `stop_fn`: A function receives the average undiscounted returns of the testing result, return a boolean which indicates whether reaching the goal.
- `writer`: See below.

The trainer supports [TensorBoard](#) for logging. It can be used as:

```
from torch.utils.tensorboard import SummaryWriter
writer = SummaryWriter('log/dqn')
```

Pass the writer into the trainer, and the training result will be recorded into the TensorBoard.

The returned result is a dictionary as follows:

```
{
    'train_step': 9246,
    'train_episode': 504.0,
    'train_time/collector': '0.65s',
    'train_time/model': '1.97s',
    'train_speed': '3518.79 step/s',
    'test_step': 49112,
    'test_episode': 400.0,
    'test_time': '1.38s',
    'test_speed': '35600.52 step/s',
    'best_reward': 199.03,
    'duration': '4.01s'
}
```

It shows that within approximately 4 seconds, we finished training a DQN agent on CartPole. The mean returns over 100 consecutive episodes is 199.03.

1.1.7 Save/Load Policy

Since the policy inherits the `torch.nn.Module` class, saving and loading the policy are exactly the same as a torch module:

```
torch.save(policy.state_dict(), 'dqn.pth')
policy.load_state_dict(torch.load('dqn.pth'))
```

1.1.8 Watch the Agent's Performance

`Collector` supports rendering. Here is the example of watching the agent's performance in 35 FPS:

```
collector = ts.data.Collector(policy, env)
collector.collect(n_episode=1, render=1 / 35)
collector.close()
```

1.1.9 Train a Policy with Customized Codes

“I don't want to use your provided trainer. I want to customize it!”

No problem! Tianshou supports user-defined training code. Here is the usage:

```
# pre-collect 5000 frames with random action before training
policy.set_eps(1)
train_collector.collect(n_step=5000)

policy.set_eps(0.1)
for i in range(int(1e6)): # total step
    collect_result = train_collector.collect(n_step=10)

    # once if the collected episodes' mean returns reach the threshold,
    # or every 1000 steps, we test it on test_collector
    if collect_result['rew'] >= env.spec.reward_threshold or i % 1000 == 0:
        policy.set_eps(0.05)
        result = test_collector.collect(n_episode=100)
        if result['rew'] >= env.spec.reward_threshold:
            print(f'Finished training! Test mean returns: {result["rew"]}')
            break
        else:
            # back to training eps
            policy.set_eps(0.1)

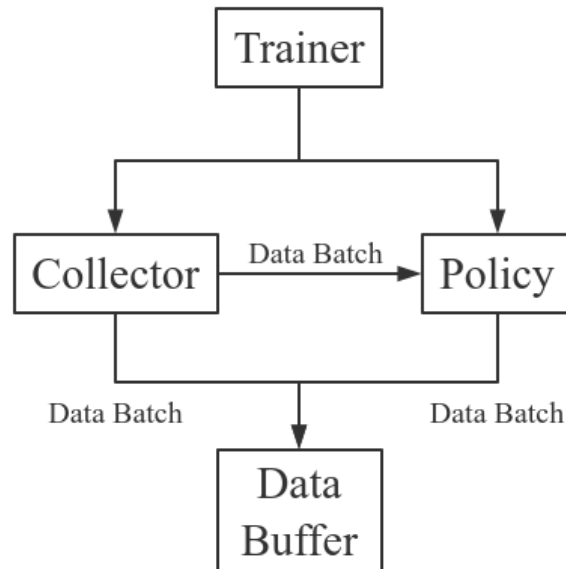
    # train policy with a sampled batch data
    losses = policy.learn(train_collector.sample(batch_size=64))
```

For further usage, you can refer to [Cheat Sheet](#).

References

1.2 Basic concepts in Tianshou

Tianshou splits a Reinforcement Learning agent training procedure into these parts: trainer, collector, policy, and data buffer. The general control flow can be described as:



1.2.1 Data Batch

Tianshou provides `Batch` as the internal data structure to pass any kind of data to other methods, for example, a collector gives a `Batch` to policy for learning. Here is the usage:

```
>>> import numpy as np
>>> from tianshou.data import Batch
>>> data = Batch(a=4, b=[5, 5], c='2312312')
>>> data.b
[5, 5]
>>> data.b = np.array([3, 4, 5])
>>> print(data)
Batch(
  a: 4,
  b: array([3, 4, 5]),
  c: '2312312',
)
```

In short, you can define a `Batch` with any key-value pair. The current implementation of Tianshou typically use 7 reserved keys in `Batch`:

- `obs` the observation of step t ;
- `act` the action of step t ;
- `rew` the reward of step t ;
- `done` the done flag of step t ;

- `obs_next` the observation of step $t + 1$;
- `info` the info of step t (in `gym.Env`, the `env.step()` function return 4 arguments, and the last one is `info`);
- `policy` the data computed by policy in step t ;

`Batch` has other methods, including `__getitem__()`, `__len__()`, `append()`, and `split()`:

```
>>> data = Batch(obs=np.array([0, 11, 22]), rew=np.array([6, 6, 6]))
>>> # here we test __getitem__
>>> index = [2, 1]
>>> data[index].obs
array([22, 11])

>>> # here we test __len__
>>> len(data)
3

>>> data.append(data) # similar to list.append
>>> data.obs
array([0, 11, 22, 0, 11, 22])

>>> # split whole data into multiple small batch
>>> for d in data.split(size=2, shuffle=False):
...     print(d.obs, d.rew)
[ 0 11] [6 6]
[22  0] [6 6]
[11 22] [6 6]
```

1.2.2 Data Buffer

`ReplayBuffer` stores data generated from interaction between the policy and environment. It stores basically 7 types of data, as mentioned in `Batch`, based on `numpy.ndarray`. Here is the usage:

```
>>> import numpy as np
>>> from tianshou.data import ReplayBuffer
>>> buf = ReplayBuffer(size=20)
>>> for i in range(3):
...     buf.add(obs=i, act=i, rew=i, done=i, obs_next=i + 1, info={})
>>> len(buf)
3
>>> buf.obs
# since we set size = 20, len(buf.obs) == 20.
array([0., 1., 2., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
       0., 0., 0., 0.])

>>> buf2 = ReplayBuffer(size=10)
>>> for i in range(15):
...     buf2.add(obs=i, act=i, rew=i, done=i, obs_next=i + 1, info={})
>>> len(buf2)
10
>>> buf2.obs
# since its size = 10, it only stores the last 10 steps' result.
array([10., 11., 12., 13., 14.,  5.,  6.,  7.,  8.,  9.])

>>> # move buf2's result into buf (meanwhile keep it chronologically)
>>> buf.update(buf2)
```

(continues on next page)

(continued from previous page)

```

array([ 0.,  1.,  2.,  5.,  6.,  7.,  8.,  9., 10., 11., 12., 13., 14.,
        0.,  0.,  0.,  0.,  0.,  0.,  0.])

>>> # get a random sample from buffer
>>> # the batch_data is equal to buf[indice].
>>> batch_data, indice = buf.sample(batch_size=4)
>>> batch_data.obs == buf[indice].obs
array([ True,  True,  True,  True])

```

ReplayBuffer also supports *frame_stack* sampling (typically for RNN usage, see issue#19), ignoring storing the next observation (save memory in atari tasks), and multi-modal observation (see issue#38, need version $\geq 0.2.3$):

```

>>> buf = ReplayBuffer(size=9, stack_num=4, ignore_obs_next=True)
>>> for i in range(16):
...     done = i % 5 == 0
...     buf.add(obs={'id': i}, act=i, rew=i, done=done,
...             obs_next={'id': i + 1})
>>> print(buf) # you can see obs_next is not saved in buf
ReplayBuffer(
  act: array([ 9., 10., 11., 12., 13., 14., 15.,  7.,  8.]),
  done: array([0.,  1.,  0.,  0.,  0.,  0.,  1.,  0.,  0.]),
  info: Batch(),
  obs: Batch(
    id: array([ 9., 10., 11., 12., 13., 14., 15.,  7.,  8.]),
  ),
  policy: Batch(),
  rew: array([ 9., 10., 11., 12., 13., 14., 15.,  7.,  8.]),
)
>>> index = np.arange(len(buf))
>>> print(buf.get(index, 'obs').id)
[[ 7.  7.  8.  9.]
 [ 7.  8.  9. 10.]
 [11. 11. 11. 11.]
 [11. 11. 11. 12.]
 [11. 11. 12. 13.]
 [11. 12. 13. 14.]
 [12. 13. 14. 15.]
 [ 7.  7.  7.  7.]
 [ 7.  7.  7.  8.]]
>>> # here is another way to get the stacked data
>>> # (stack only for obs and obs_next)
>>> abs(buf.get(index, 'obs')['id'] - buf[index].obs.id).sum().sum()
0.0
>>> # we can get obs_next through __getitem__, even if it doesn't exist
>>> print(buf[:].obs_next.id)
[[ 7.  8.  9. 10.]
 [ 7.  8.  9. 10.]
 [11. 11. 11. 12.]
 [11. 11. 12. 13.]
 [11. 12. 13. 14.]
 [12. 13. 14. 15.]
 [12. 13. 14. 15.]
 [ 7.  7.  7.  8.]
 [ 7.  7.  8.  9.]]

```

Tianshou provides other type of data buffer such as *ListReplayBuffer* (based on list), *PrioritizedReplayBuffer* (based on Segment Tree and `numpy.ndarray`). Check out *ReplayBuffer*

for more detail.

1.2.3 Policy

Tianshou aims to modularizing RL algorithms. It comes into several classes of policies in Tianshou. All of the policy classes must inherit `BasePolicy`.

A policy class typically has four parts:

- `__init__()`: initialize the policy, including coping the target network and so on;
- `forward()`: compute action with given observation;
- `process_fn()`: pre-process data from the replay buffer (this function can interact with replay buffer);
- `learn()`: update policy with a given batch of data.

Take 2-step return DQN as an example. The 2-step return DQN compute each frame's return as:

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 \max_a Q(s_{t+2}, a)$$

where γ is the discount factor, $\gamma \in [0, 1]$. Here is the pseudocode showing the training process **without Tianshou framework**:

```
# pseudocode, cannot work
s = env.reset()
buffer = Buffer(size=10000)
agent = DQN()
for i in range(int(1e6)):
    a = agent.compute_action(s)
    s_, r, d, _ = env.step(a)
    buffer.store(s, a, s_, r, d)
    s = s_
    if i % 1000 == 0:
        b_s, b_a, b_s_, b_r, b_d = buffer.get(size=64)
        # compute 2-step returns. How?
        b_ret = compute_2_step_return(buffer, b_r, b_d, ...)
        # update DQN policy
        agent.update(b_s, b_a, b_s_, b_r, b_d, b_ret)
```

Thus, we need a time-related interface for calculating the 2-step return. `process_fn()` finishes this work by providing the replay buffer, the sample index, and the sample batch data. Since we store all the data in the order of time, you can simply compute the 2-step return as:

```
class DQN_2step(BasePolicy):
    """some code"""

    def process_fn(self, batch, buffer, indice):
        buffer_len = len(buffer)
        batch_2 = buffer[(indice + 2) % buffer_len]
        # this will return a batch data where batch_2.obs is s_t+2
        # we can also get s_t+2 through:
        # batch_2_obs = buffer.obs[(indice + 2) % buffer_len]
        # in short, buffer.obs[i] is equal to buffer[i].obs, but the former is more_
        ↪ efficient.
        Q = self(batch_2, eps=0) # shape: [batchsize, action_shape]
        maxQ = Q.max(dim=-1)
        batch.returns = batch.rew \
```

(continues on next page)

(continued from previous page)

```

        + self._gamma * buffer.rew[(indice + 1) % buffer_len] \
        + self._gamma ** 2 * maxQ
    return batch

```

This code does not consider the done flag, so it may not work very well. It shows two ways to get s_{t+2} from the replay buffer easily in `process_fn()`.

For other method, you can check out [tianshou.policy](#). We give the usage of policy class a high-level explanation in [A High-level Explanation](#).

1.2.4 Collector

The `Collector` enables the policy to interact with different types of environments conveniently. In short, `Collector` has two main methods:

- `collect()`: let the policy perform (at least) a specified number of step `n_step` or episode `n_episode` and store the data in the replay buffer;
- `sample()`: sample a data batch from replay buffer; it will call `process_fn()` before returning the final batch data.

Why do we mention **at least** here? For a single environment, the collector will finish exactly `n_step` or `n_episode`. However, for multiple environments, we could not directly store the collected data into the replay buffer, since it breaks the principle of storing data chronologically.

The solution is to add some cache buffers inside the collector. Once collecting **a full episode of trajectory**, it will move the stored data from the cache buffer to the main buffer. To satisfy this condition, the collector will interact with environments that may exceed the given step number or episode number.

The general explanation is listed in [A High-level Explanation](#). Other usages of collector are listed in `Collector` documentation.

1.2.5 Trainer

Once you have a collector and a policy, you can start writing the training method for your RL agent. Trainer, to be honest, is a simple wrapper. It helps you save energy for writing the training loop. You can also construct your own trainer: [Train a Policy with Customized Codes](#).

Tianshou has two types of trainer: `onpolicy_trainer()` and `offpolicy_trainer()`, corresponding to on-policy algorithms (such as Policy Gradient) and off-policy algorithms (such as DQN). Please check out [tianshou.trainer](#) for the usage.

There will be more types of trainers, for instance, multi-agent trainer.

1.2.6 A High-level Explanation

We give a high-level explanation through the pseudocode used in section [Policy](#):

```

# pseudocode, cannot work                                     # methods in tianshou
s = env.reset()                                                # buffer = tianshou.
buffer = Buffer(size=10000)
↪ data.ReplayBuffer(size=10000)
agent = DQN()                                                  # policy.__init__(...)
for i in range(int(1e6)):                                       # done in trainer

```

(continues on next page)

(continued from previous page)

```

a = agent.compute_action(s)                # policy(batch, ...)
s_, r, d, _ = env.step(a)                  # collector.collect(..
↪.)
buffer.store(s, a, s_, r, d)                # collector.collect(..
↪.)
s = s_                                     # collector.collect(..
↪.)
    if i % 1000 == 0:                       # done in trainer
        b_s, b_a, b_s_, b_r, b_d = buffer.get(size=64) # collector.
↪sample(batch_size)
        # compute 2-step returns. How?
        b_ret = compute_2_step_return(buffer, b_r, b_d, ...) # policy.process_
↪fn(batch, buffer, indice)
        # update DQN policy
        agent.update(b_s, b_a, b_s_, b_r, b_d, b_ret)      # policy.learn(batch,
↪...)

```

1.2.7 Conclusion

So far, we go through the overall framework of Tianshou. Really simple, isn't it?

1.3 Train a model-free RL agent within 30s

This page summarizes some hyper-parameter tuning experience and code-level trick when training a model-free DRL agent.

You can also contribute to this page with your own tricks :)

1.3.1 Avoid batch-size = 1

In the traditional RL training loop, we always use the policy to interact with only one environment for collecting data. That means most of the time the network use batch-size = 1. Quite inefficient! Here is an example of showing how inefficient it is:

```

import torch, time
from torch import nn

class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.model = nn.Sequential(
            nn.Linear(3, 128), nn.ReLU(inplace=True),
            nn.Linear(128, 128), nn.ReLU(inplace=True),
            nn.Linear(128, 1))
    def forward(self, s):
        return self.model(s)

net = Net()
cnt = 1000
div = 128
a = torch.randn([128, 3])

```

(continues on next page)

(continued from previous page)

```

t = time.time()
for i in range(cnt):
    b = net(a)
t1 = (time.time() - t) / cnt
print(t1)
t = time.time()
for i in range(cnt):
    for a_ in a.split(a.shape[0] // div):
        b = net(a_)
t2 = (time.time() - t) / cnt
print(t2)
print(t2 / t1)

```

The first test uses batch-size 128, and the second test uses batch-size = 1 for 128 times. In our test, the first is 70-80 times faster than the second.

So how could we avoid the case of batch-size = 1? The answer is synchronize sampling: we create multiple independent environments and sample simultaneously. It is similar to A2C, but other algorithms can also use this method. In our experiments, sampling from more environments benefits not only the sample speed but also the converge speed of neural network (we guess it lowers the sample bias).

By the way, A2C is better than A3C in some cases: A3C needs to act independently and sync the gradient to master, but, in a single node, using A3C to act with batch-size = 1 is quite resource-consuming.

1.3.2 Algorithm specific tricks

Here is about the experience of hyper-parameter tuning on CartPole and Pendulum:

- *DQNPolicy*: use estimation_step greater than 1 and target network, also with a suitable size of replay buffer;
- *PGPolicy*: TBD
- *A2CPolicy*: TBD
- *PPOPolicy*: TBD
- *DDPGPolicy*, *TD3Policy*, and *SACPolicy*: We found two tricks. The first is to ignore the done flag. The second is to normalize reward to a standard normal distribution (it is against the theoretical analysis, but indeed works very well). The two tricks work amazingly on Mujoco tasks, typically with a faster converge speed (1M -> 200K).
- On-policy algorithms: increase the repeat-time (to 2 or 4 for trivial benchmark, 10 for mujoco) of the given batch in each training update will make the algorithm more stable.

1.3.3 Code-level optimization

Tianshou has many short-but-efficient lines of code. For example, when we want to compute $V(s)$ and $V(s')$ by the same network, the best way is to concatenate s and s' together instead of computing the value function using twice of network forward.

1.3.4 Finally

With fast-speed sampling, we could use large batch-size and large learning rate for faster convergence.

RL algorithms are seed-sensitive. Try more seeds and pick the best. But for our demo, we just used seed = 0 and found it work surprisingly well on policy gradient, so we did not try other seed.

1.4 Cheat Sheet

This page shows some code snippets of how to use Tianshou to develop new algorithms.

TODO

1.5 tianshou.data

class tianshou.data.Batch (**kwargs)

Bases: object

Tianshou provides *Batch* as the internal data structure to pass any kind of data to other methods, for example, a collector gives a *Batch* to policy for learning. Here is the usage:

```
>>> import numpy as np
>>> from tianshou.data import Batch
>>> data = Batch(a=4, b=[5, 5], c='2312312')
>>> data.b
[5, 5]
>>> data.b = np.array([3, 4, 5])
>>> print(data)
Batch(
  a: 4,
  b: array([3, 4, 5]),
  c: '2312312',
)
```

In short, you can define a *Batch* with any key-value pair. The current implementation of Tianshou typically use 7 reserved keys in *Batch*:

- *obs* the observation of step t ;
- *act* the action of step t ;
- *rew* the reward of step t ;
- *done* the done flag of step t ;
- *obs_next* the observation of step $t + 1$;
- *info* the info of step t (in `gym.Env`, the `env.step()` function return 4 arguments, and the last one is `info`);
- *policy* the data computed by policy in step t ;

Batch has other methods, including `__getitem__()`, `__len__()`, `append()`, and `split()`:

```

>>> data = Batch(obs=np.array([0, 11, 22]), rew=np.array([6, 6, 6]))
>>> # here we test __getitem__
>>> index = [2, 1]
>>> data[index].obs
array([22, 11])

>>> # here we test __len__
>>> len(data)
3

>>> data.append(data) # similar to list.append
>>> data.obs
array([0, 11, 22, 0, 11, 22])

>>> # split whole data into multiple small batch
>>> for d in data.split(size=2, shuffle=False):
...     print(d.obs, d.rew)
[ 0 11] [6 6]
[22  0] [6 6]
[11 22] [6 6]

```

__getitem__ (*index: Union[str, slice]*) → Union[tianshou.data.batch.Batch, dict]
Return self[index].

__len__ () → int
Return len(self).

append (*batch: tianshou.data.batch.Batch*) → None
Append a *Batch* object to current batch.

get (*k: str, d: Optional[Any] = None*) → Union[tianshou.data.batch.Batch, Any]
Return self[k] if k in self else d. d defaults to None.

keys () → List[str]
Return self.keys().

split (*size: Optional[int] = None, shuffle: bool = True*) → Iterator[tianshou.data.batch.Batch]
Split whole data into multiple small batch.

Parameters

- **size** (*int*) – if it is None, it does not split the data batch; otherwise it will divide the data batch with the given size. Default to None.
- **shuffle** (*bool*) – randomly shuffle the entire data batch if it is True, otherwise remain in the same. Default to True.

to_numpy () → None
Change all torch.Tensor to numpy.ndarray. This is an inplace operation.

to_torch (*dtype: Optional[torch.dtype] = None, device: Union[str, int, torch.device] = 'cpu'*) → None
Change all numpy.ndarray to torch.Tensor. This is an inplace operation.

values () → List[Any]
Return self.values().

tianshou.data.to_numpy (*x: Union[torch.Tensor, dict, tianshou.data.batch.Batch, numpy.ndarray]*) → Union[dict, tianshou.data.batch.Batch, numpy.ndarray]
Return an object without torch.Tensor.

`tianshou.data.to_torch` (*x*: `Union[torch.Tensor, dict, tianshou.data.batch.Batch, numpy.ndarray]`,
dtype: `Optional[torch.dtype] = None`, *device*: `Union[str, int] = 'cpu'`) →
`Union[dict, tianshou.data.batch.Batch, torch.Tensor]`

Return an object without `np.ndarray`.

class `tianshou.data.ReplayBuffer` (*size*: `int`, *stack_num*: `Optional[int] = 0`, *ignore_obs_next*:
bool = `False`, ***kwargs*)

Bases: `object`

`ReplayBuffer` stores data generated from interaction between the policy and environment. It stores basically 7 types of data, as mentioned in `Batch`, based on `numpy.ndarray`. Here is the usage:

```
>>> import numpy as np
>>> from tianshou.data import ReplayBuffer
>>> buf = ReplayBuffer(size=20)
>>> for i in range(3):
...     buf.add(obs=i, act=i, rew=i, done=i, obs_next=i + 1, info={})
>>> len(buf)
3
>>> buf.obs
# since we set size = 20, len(buf.obs) == 20.
array([0., 1., 2., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,
       0., 0., 0., 0.])

>>> buf2 = ReplayBuffer(size=10)
>>> for i in range(15):
...     buf2.add(obs=i, act=i, rew=i, done=i, obs_next=i + 1, info={})
>>> len(buf2)
10
>>> buf2.obs
# since its size = 10, it only stores the last 10 steps' result.
array([10., 11., 12., 13., 14., 5., 6., 7., 8., 9.])

>>> # move buf2's result into buf (meanwhile keep it chronologically)
>>> buf.update(buf2)
array([ 0., 1., 2., 5., 6., 7., 8., 9., 10., 11., 12., 13., 14.,
       0., 0., 0., 0., 0., 0., 0.])

>>> # get a random sample from buffer
>>> # the batch_data is equal to buf[indice].
>>> batch_data, indice = buf.sample(batch_size=4)
>>> batch_data.obs == buf[indice].obs
array([ True,  True,  True,  True])
```

`ReplayBuffer` also supports `frame_stack` sampling (typically for RNN usage, see issue#19), ignoring storing the next observation (save memory in atari tasks), and multi-modal observation (see issue#38, need version >= 0.2.3):

```
>>> buf = ReplayBuffer(size=9, stack_num=4, ignore_obs_next=True)
>>> for i in range(16):
...     done = i % 5 == 0
...     buf.add(obs={'id': i}, act=i, rew=i, done=done,
...             obs_next={'id': i + 1})
>>> print(buf) # you can see obs_next is not saved in buf
ReplayBuffer(
  act: array([ 9., 10., 11., 12., 13., 14., 15., 7., 8.]),
  done: array([0., 1., 0., 0., 0., 0., 1., 0., 0.]),
  info: Batch(),
```

(continues on next page)

(continued from previous page)

```

obs: Batch(
    id: array([ 9., 10., 11., 12., 13., 14., 15.,  7.,  8.]),
),
policy: Batch(),
rew: array([ 9., 10., 11., 12., 13., 14., 15.,  7.,  8.]),
)
>>> index = np.arange(len(buf))
>>> print(buf.get(index, 'obs').id)
[[ 7.  7.  8.  9.]
 [ 7.  8.  9. 10.]
 [11. 11. 11. 11.]
 [11. 11. 11. 12.]
 [11. 11. 12. 13.]
 [11. 12. 13. 14.]
 [12. 13. 14. 15.]
 [ 7.  7.  7.  7.]
 [ 7.  7.  7.  8.]]
>>> # here is another way to get the stacked data
>>> # (stack only for obs and obs_next)
>>> abs(buf.get(index, 'obs')['id'] - buf[index].obs.id).sum().sum()
0.0
>>> # we can get obs_next through __getitem__, even if it doesn't exist
>>> print(buf[:].obs_next.id)
[[ 7.  8.  9. 10.]
 [ 7.  8.  9. 10.]
 [11. 11. 11. 12.]
 [11. 11. 12. 13.]
 [11. 12. 13. 14.]
 [12. 13. 14. 15.]
 [12. 13. 14. 15.]
 [ 7.  7.  7.  8.]
 [ 7.  7.  8.  9.]]

```

__getitem__ (*index: Union[slice, numpy.ndarray]*) → *tianshou.data.batch.Batch*

Return a data batch: `self[index]`. If `stack_num` is set to be `> 0`, return the stacked `obs` and `obs_next` with shape `[batch, len, ...]`.

__len__ () → *int*

Return `len(self)`.

add (*obs: Union[dict, numpy.ndarray]*, *act: Union[numpy.ndarray, float]*, *rew: float*, *done: bool*, *obs_next: Union[dict, numpy.ndarray, None] = None*, *info: dict = {}*, *policy: Union[dict, tianshou.data.batch.Batch, None] = {}*, ***kwargs*) → *None*
Add a batch of data into replay buffer.

get (*indice: Union[slice, numpy.ndarray]*, *key: str*, *stack_num: Optional[int] = None*) → *Union[tianshou.data.batch.Batch, numpy.ndarray]*

Return the stacked result, e.g. `[s_{t-3}, s_{t-2}, s_{t-1}, s_t]`, where `s` is `self.key`, `t` is `indice`. The `stack_num` (here equals to 4) is given from buffer initialization procedure.

reset () → *None*

Clear all the data in replay buffer.

sample (*batch_size: int*) → *Tuple[tianshou.data.batch.Batch, numpy.ndarray]*

Get a random sample from buffer with size equal to `batch_size`. Return all the data in the buffer if `batch_size` is 0.

Returns Sample data and its corresponding index inside the buffer.

update (*buffer: tianshou.data.buffer.ReplayBuffer*) → None
Move the data from the given buffer to self.

class `tianshou.data.ListReplayBuffer` (***kwargs*)
Bases: `tianshou.data.buffer.ReplayBuffer`

The function of `ListReplayBuffer` is almost the same as `ReplayBuffer`. The only difference is that `ListReplayBuffer` is based on list.

See also:

Please refer to `ReplayBuffer` for more detailed explanation.

reset () → None
Clear all the data in replay buffer.

class `tianshou.data.PrioritizedReplayBuffer` (*size: int, alpha: float, beta: float, mode: str = 'weight', **kwargs*)
Bases: `tianshou.data.buffer.ReplayBuffer`

Prioritized replay buffer implementation.

Parameters

- **alpha** (*float*) – the prioritization exponent.
- **beta** (*float*) – the importance sample soft coefficient.
- **mode** (*str*) – defaults to weight.

See also:

Please refer to `ReplayBuffer` for more detailed explanation.

__getitem__ (*index: Union[slice, numpy.ndarray]*) → `tianshou.data.batch.Batch`
Return a data batch: `self[index]`. If `stack_num` is set to be > 0, return the stacked obs and obs_next with shape `[batch, len, ...]`.

add (*obs: Union[dict, numpy.ndarray], act: Union[numpy.ndarray, float], rew: float, done: bool, obs_next: Union[dict, numpy.ndarray, None] = None, info: dict = {}, policy: Union[dict, tianshou.data.batch.Batch, None] = {}, weight: float = 1.0, **kwargs*) → None
Add a batch of data into replay buffer.

reset () → None
Clear all the data in replay buffer.

sample (*batch_size: int, importance_sample: bool = True*) → `Tuple[tianshou.data.batch.Batch, numpy.ndarray]`
Get a random sample from buffer with priority probability. Return all the data in the buffer if `batch_size` is 0.

Returns Sample data and its corresponding index inside the buffer.

update_weight (*indice: Union[slice, numpy.ndarray], new_weight: numpy.ndarray*) → None
Update priority weight by indice in this buffer.

Parameters

- **indice** (*np.ndarray*) – indice you want to update weight
- **new_weight** (*np.ndarray*) – new priority weight you want to update

```
class tianshou.data.Collector(policy:          tianshou.policy.base.BasePolicy,          env:
                             Union[gym.core.Env,          tianshou.env.vecenv.BaseVectorEnv],
                             buffer:          Union[tianshou.data.buffer.ReplayBuffer,
                             List[tianshou.data.buffer.ReplayBuffer], None] = None, preprocess_fn: Callable[[Any], Union[dict, tianshou.data.batch.Batch]]
                             = None, stat_size: Optional[int] = 100, **kwargs)
```

Bases: object

The `Collector` enables the policy to interact with different types of environments conveniently.

Parameters

- **policy** – an instance of the `BasePolicy` class.
- **env** – a `gym.Env` environment or an instance of the `BaseVectorEnv` class.
- **buffer** – an instance of the `ReplayBuffer` class, or a list of `ReplayBuffer`. If set to `None`, it will automatically assign a small-size `ReplayBuffer`.
- **preprocess_fn** (*function*) – a function called before the data has been added to the buffer, see issue #42, defaults to `None`.
- **stat_size** (*int*) – for the moving average of recording speed, defaults to 100.

The `preprocess_fn` is a function called before the data has been added to the buffer with batch format, which receives up to 7 keys as listed in `Batch`. It will receive with only `obs` when the collector resets the environment. It returns either a dict or a `Batch` with the modified keys and values. Examples are in “test/base/test_collector.py”.

Example:

```
policy = PGPolicy(...) # or other policies if you wish
env = gym.make('CartPole-v0')
replay_buffer = ReplayBuffer(size=10000)
# here we set up a collector with a single environment
collector = Collector(policy, env, buffer=replay_buffer)

# the collector supports vectorized environments as well
envs = VectorEnv([lambda: gym.make('CartPole-v0') for _ in range(3)])
buffers = [ReplayBuffer(size=5000) for _ in range(3)]
# you can also pass a list of replay buffer to collector, for multi-env
# collector = Collector(policy, envs, buffer=buffers)
collector = Collector(policy, envs, buffer=replay_buffer)

# collect at least 3 episodes
collector.collect(n_episode=3)
# collect 1 episode for the first env, 3 for the third env
collector.collect(n_episode=[1, 0, 3])
# collect at least 2 steps
collector.collect(n_step=2)
# collect episodes with visual rendering (the render argument is the
#   sleep time between rendering consecutive frames)
collector.collect(n_episode=1, render=0.03)

# sample data with a given number of batch-size:
batch_data = collector.sample(batch_size=64)
# policy.learn(batch_data) # btw, vanilla policy gradient only
#   supports on-policy training, so here we pick all data in the buffer
batch_data = collector.sample(batch_size=0)
policy.learn(batch_data)
```

(continues on next page)

(continued from previous page)

```
# on-policy algorithms use the collected data only once, so here we
# clear the buffer
collector.reset_buffer()
```

For the scenario of collecting data from multiple environments to a single buffer, the cache buffers will turn on automatically. It may return the data more than the given limitation.

Note: Please make sure the given environment has a time limitation.

close() → None

Close the environment(s).

collect (*n_step*: int = 0, *n_episode*: Union[int, List[int]] = 0, *render*: Optional[float] = None, *log_fn*: Optional[Callable[[dict], None]] = None) → Dict[str, float]

Collect a specified number of step or episode.

Parameters

- **n_step** (*int*) – how many steps you want to collect.
- **n_episode** (*int or list*) – how many episodes you want to collect (in each environment).
- **render** (*float*) – the sleep time between rendering consecutive frames, defaults to None (no rendering).
- **log_fn** (*function*) – a function which receives env info, typically for tensorboard logging.

Note: One and only one collection number specification is permitted, either *n_step* or *n_episode*.

Returns

A dict including the following keys

- *n/ep* the collected number of episodes.
- *n/st* the collected number of steps.
- *v/st* the speed of steps per second.
- *v/ep* the speed of episode per second.
- *rew* the mean reward over collected episodes.
- *len* the mean length over collected episodes.

get_env_num() → int

Return the number of environments the collector have.

render (***kwargs*) → None

Render all the environment(s).

reset() → None

Reset all related variables in the collector.

reset_buffer() → None

Reset the main data buffer.

reset_env () → None

Reset all of the environment(s)' states and reset all of the cache buffers (if need).

sample (batch_size: int) → tianshou.data.batch.Batch

Sample a data batch from the internal replay buffer. It will call `process_fn` () before returning the final batch data.

Parameters **batch_size** (int) - 0 means it will extract all the data from the buffer, otherwise it will extract the data with the given batch_size.

seed (seed: Union[int, List[int], None] = None) → None

Reset all the seed(s) of the given environment(s).

1.6 tianshou.env

class tianshou.env.BaseVectorEnv (env_fns: List[Callable[], gym.core.Env])

Bases: abc.ABC, gym.core.Env

Base class for vectorized environments wrapper. Usage:

```
env_num = 8
envs = VectorEnv([lambda: gym.make(task) for _ in range(env_num)])
assert len(envs) == env_num
```

It accepts a list of environment generators. In other words, an environment generator `efn` of a specific task means that `efn()` returns the environment of the given task, for example, `gym.make(task)`.

All of the VectorEnv must inherit `BaseVectorEnv`. Here are some other usages:

```
envs.seed(2) # which is equal to the next line
envs.seed([2, 3, 4, 5, 6, 7, 8, 9]) # set specific seed for each env
obs = envs.reset() # reset all environments
obs = envs.reset([0, 5, 7]) # reset 3 specific environments
obs, rew, done, info = envs.step([1] * 8) # step synchronously
envs.render() # render all environments
envs.close() # close all environments
```

__len__ () → int

Return len(self), which is the number of environments.

abstract close () → None

Close all of the environments.

Environments will automatically close() themselves when garbage collected or when the program exits.

abstract render (**kwargs) → None

Render all of the environments.

abstract reset (id: Union[int, List[int], None] = None)

Reset the state of all the environments and return initial observations if id is None, otherwise reset the specific environments with given id, either an int or a list.

abstract seed (seed: Union[int, List[int], None] = None) → None

Set the seed for all environments.

Accept None, an int (which will extend i to [i, i + 1, i + 2, ...]) or a list.

Returns The list of seeds used in this env's random number generators.

The first value in the list should be the "main" seed, or the value which a reproducer should pass to 'seed'.

abstract step (*action*: *numpy.ndarray*) → Tuple[*numpy.ndarray*, *numpy.ndarray*, *numpy.ndarray*, *numpy.ndarray*]

Run one timestep of all the environments' dynamics. When the end of episode is reached, you are responsible for calling `reset(id)` to reset this environment's state.

Accept a batch of action and return a tuple (obs, rew, done, info).

Parameters *action* (*numpy.ndarray*) – a batch of action provided by the agent.

Returns

A tuple including four items:

- *obs* a *numpy.ndarray*, the agent's observation of current environments
- *rew* a *numpy.ndarray*, the amount of rewards returned after previous actions
- *done* a *numpy.ndarray*, whether these episodes have ended, in which case further `step()` calls will return undefined results
- *info* a *numpy.ndarray*, contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

class `tianshou.env.VectorEnv` (*env_fns*: List[Callable[], *gym.core.Env*])

Bases: `tianshou.env.vecenv.BaseVectorEnv`

Dummy vectorized environment wrapper, implemented in for-loop.

See also:

Please refer to [BaseVectorEnv](#) for more detailed explanation.

close () → None

Close all of the environments.

Environments will automatically `close()` themselves when garbage collected or when the program exits.

render (***kwargs*) → None

Render all of the environments.

reset (*id*: Union[int, List[int], None] = None) → None

Reset the state of all the environments and return initial observations if *id* is None, otherwise reset the specific environments with given *id*, either an int or a list.

seed (*seed*: Union[int, List[int], None] = None) → None

Set the seed for all environments.

Accept None, an int (which will extend *i* to [*i*, *i* + 1, *i* + 2, ...]) or a list.

Returns The list of seeds used in this env's random number generators.

The first value in the list should be the “main” seed, or the value which a reproducer should pass to ‘seed’.

step (*action*: *numpy.ndarray*) → Tuple[*numpy.ndarray*, *numpy.ndarray*, *numpy.ndarray*, *numpy.ndarray*]

Run one timestep of all the environments' dynamics. When the end of episode is reached, you are responsible for calling `reset(id)` to reset this environment's state.

Accept a batch of action and return a tuple (obs, rew, done, info).

Parameters *action* (*numpy.ndarray*) – a batch of action provided by the agent.

Returns

A tuple including four items:

- *obs* a *numpy.ndarray*, the agent's observation of current environments

- `rew` a `numpy.ndarray`, the amount of rewards returned after previous actions
- `done` a `numpy.ndarray`, whether these episodes have ended, in which case further `step()` calls will return undefined results
- `info` a `numpy.ndarray`, contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

class `tianshou.env.SubprocVectorEnv` (*env_fns: List[Callable[], gym.core.Env[]]*)

Bases: `tianshou.env.vecenv.BaseVectorEnv`

Vectorized environment wrapper based on subprocess.

See also:

Please refer to [`BaseVectorEnv`](#) for more detailed explanation.

close () → None

Close all of the environments.

Environments will automatically `close()` themselves when garbage collected or when the program exits.

render (***kwargs*) → None

Render all of the environments.

reset (*id: Union[int, List[int], None] = None*) → None

Reset the state of all the environments and return initial observations if `id` is `None`, otherwise reset the specific environments with given `id`, either an `int` or a list.

seed (*seed: Union[int, List[int], None] = None*) → None

Set the seed for all environments.

Accept `None`, an `int` (which will extend `i` to `[i, i + 1, i + 2, ...]`) or a list.

Returns The list of seeds used in this env’s random number generators.

The first value in the list should be the “main” seed, or the value which a reproducer should pass to ‘seed’.

step (*action: numpy.ndarray*) → Tuple[`numpy.ndarray`, `numpy.ndarray`, `numpy.ndarray`, `numpy.ndarray`]

Run one timestep of all the environments’ dynamics. When the end of episode is reached, you are responsible for calling `reset(id)` to reset this environment’s state.

Accept a batch of action and return a tuple (`obs`, `rew`, `done`, `info`).

Parameters **action** (*numpy.ndarray*) – a batch of action provided by the agent.

Returns

A tuple including four items:

- `obs` a `numpy.ndarray`, the agent’s observation of current environments
- `rew` a `numpy.ndarray`, the amount of rewards returned after previous actions
- `done` a `numpy.ndarray`, whether these episodes have ended, in which case further `step()` calls will return undefined results
- `info` a `numpy.ndarray`, contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

class `tianshou.env.RayVectorEnv` (*env_fns: List[Callable[], gym.core.Env[]]*)

Bases: `tianshou.env.vecenv.BaseVectorEnv`

Vectorized environment wrapper based on [`ray`](#). However, according to our test, it is about two times slower than [`SubprocVectorEnv`](#).

See also:

Please refer to [BaseVectorEnv](#) for more detailed explanation.

close () → None

Close all of the environments.

Environments will automatically close() themselves when garbage collected or when the program exits.

render (**kwargs) → None

Render all of the environments.

reset (id: Union[int, List[int], None] = None) → None

Reset the state of all the environments and return initial observations if id is None, otherwise reset the specific environments with given id, either an int or a list.

seed (seed: Union[int, List[int], None] = None) → None

Set the seed for all environments.

Accept None, an int (which will extend i to [i, i + 1, i + 2, ...]) or a list.

Returns The list of seeds used in this env’s random number generators.

The first value in the list should be the “main” seed, or the value which a reproducer should pass to ‘seed’.

step (action: numpy.ndarray) → Tuple[numpy.ndarray, numpy.ndarray, numpy.ndarray, numpy.ndarray]

Run one timestep of all the environments’ dynamics. When the end of episode is reached, you are responsible for calling reset(id) to reset this environment’s state.

Accept a batch of action and return a tuple (obs, rew, done, info).

Parameters **action** (numpy.ndarray) – a batch of action provided by the agent.

Returns

A tuple including four items:

- **obs** a numpy.ndarray, the agent’s observation of current environments
- **rew** a numpy.ndarray, the amount of rewards returned after previous actions
- **done** a numpy.ndarray, whether these episodes have ended, in which case further step() calls will return undefined results
- **info** a numpy.ndarray, contains auxiliary diagnostic information (helpful for debugging, and sometimes learning)

1.7 tianshou.policy

class tianshou.policy.BasePolicy (**kwargs)

Bases: abc.ABC, torch.nn.modules.module.Module

Tianshou aims to modularizing RL algorithms. It comes into several classes of policies in Tianshou. All of the policy classes must inherit [BasePolicy](#).

A policy class typically has four parts:

- **__init__**(): initialize the policy, including coping the target network and so on;
- **forward**(): compute action with given observation;
- **process_fn**(): pre-process data from the replay buffer (this function can interact with replay buffer);

- `learn()`: update policy with a given batch of data.

Most of the policy needs a neural network to predict the action and an optimizer to optimize the policy. The rules of self-defined networks are:

1. Input: observation `obs` (may be a `numpy.ndarray` or `torch.Tensor`), hidden state `state` (for RNN usage), and other information `info` provided by the environment.
2. Output: some logits and the next hidden state `state`. The logits could be a tuple instead of a `torch.Tensor`. It depends on how the policy process the network output. For example, in PPO, the return of the network might be `(mu, sigma), state` for Gaussian policy.

Since `BasePolicy` inherits `torch.nn.Module`, you can use `BasePolicy` almost the same as `torch.nn.Module`, for instance, loading and saving the model:

```
torch.save(policy.state_dict(), 'policy.pth')
policy.load_state_dict(torch.load('policy.pth'))
```

```
static compute_episodic_return (batch: tianshou.data.batch.Batch, v_s:
                                Union[numpy.ndarray, torch.Tensor, None] = None,
                                gamma: float = 0.99, gae_lambda: float = 0.95) →
                                tianshou.data.batch.Batch
```

Compute returns over given full-length episodes, including the implementation of Generalized Advantage Estimator (arXiv:1506.02438).

Parameters

- **batch** (*Batch*) – a data batch which contains several full-episode data chronologically.
- **v_s** (*numpy.ndarray*) – the value function of all next states $V(s')$.
- **gamma** (*float*) – the discount factor, should be in $[0, 1]$, defaults to 0.99.
- **gae_lambda** (*float*) – the parameter for Generalized Advantage Estimation, should be in $[0, 1]$, defaults to 0.95.

```
abstract forward (batch: tianshou.data.batch.Batch, state: Union[dict, tianshou.data.batch.Batch,
                                                                numpy.ndarray, None] = None, **kwargs) → tianshou.data.batch.Batch
```

Compute action over the given batch data.

Returns

A *Batch* which MUST have the following keys:

- **act** an `numpy.ndarray` or a `torch.Tensor`, the action over given batch data.
- **state** a dict, an `numpy.ndarray` or a `torch.Tensor`, the internal state of the policy, `None` as default.

Other keys are user-defined. It depends on the algorithm. For example,

```
# some code
return Batch(logits=..., act=..., state=None, dist=...)
```

After version $\geq 0.2.3$, the keyword “policy” is reserved and the corresponding data will be stored into the replay buffer in `numpy`. For instance,

```
# some code
return Batch(..., policy=Batch(log_prob=dist.log_prob(act)))
# and in the sampled data batch, you can directly call
# batch.policy.log_prob to get your data, although it is stored in
# np.ndarray.
```

abstract learn (*batch*: *tianshou.data.batch.Batch*, ***kwargs*) → Dict[str, Union[float, List[float]]]
Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

process_fn (*batch*: *tianshou.data.batch.Batch*, *buffer*: *tianshou.data.buffer.ReplayBuffer*, *indice*: *numpy.ndarray*) → *tianshou.data.batch.Batch*

Pre-process the data from the provided replay buffer. Check out [Policy](#) for more information.

class *tianshou.policy.ImitationPolicy* (*model*: *torch.nn.modules.module.Module*, *optim*: *torch.optim.optimizer.Optimizer*, *mode*: *str* = 'continuous', ***kwargs*)

Bases: *tianshou.policy.base.BasePolicy*

Implementation of vanilla imitation learning (for continuous action space).

Parameters

- **model** (*torch.nn.Module*) – a model following the rules in [BasePolicy](#). (s → a)
- **optim** (*torch.optim.Optimizer*) – for optimizing the model.
- **mode** (*str*) – indicate the imitation type (“continuous” or “discrete” action space), defaults to “continuous”.

See also:

Please refer to [BasePolicy](#) for more detailed explanation.

forward (*batch*: *tianshou.data.batch.Batch*, *state*: Union[dict, *tianshou.data.batch.Batch*, *numpy.ndarray*, None] = None, ***kwargs*) → *tianshou.data.batch.Batch*
Compute action over the given batch data.

Returns

A *Batch* which MUST have the following keys:

- **act** an *numpy.ndarray* or a *torch.Tensor*, the action over given batch data.
- **state** a dict, an *numpy.ndarray* or a *torch.Tensor*, the internal state of the policy, None as default.

Other keys are user-defined. It depends on the algorithm. For example,

```
# some code
return Batch(logits=..., act=..., state=None, dist=...)
```

After version >= 0.2.3, the keyword “policy” is reserved and the corresponding data will be stored into the replay buffer in *numpy*. For instance,

```
# some code
return Batch(..., policy=Batch(log_prob=dist.log_prob(act)))
# and in the sampled data batch, you can directly call
# batch.policy.log_prob to get your data, although it is stored in
# np.ndarray.
```

learn (*batch*: *tianshou.data.batch.Batch*, ***kwargs*) → Dict[str, float]
Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

```
class tianshou.policy.DQNPolicy(model: torch.nn.modules.module.Module, optim:
    torch.optim.optimizer.Optimizer, discount_factor: float = 0.99,
    estimation_step: int = 1, target_update_freq: Optional[int] =
    0, **kwargs)
Bases: tianshou.policy.base.BasePolicy
Implementation of Deep Q Network. arXiv:1312.5602 Implementation of Double Q-Learning.
arXiv:1509.06461
```

Parameters

- **model** (*torch.nn.Module*) – a model following the rules in *BasePolicy*. (s -> logits)
- **optim** (*torch.optim.Optimizer*) – a torch.optim for optimizing the model.
- **discount_factor** (*float*) – in [0, 1].
- **estimation_step** (*int*) – greater than 1, the number of steps to look ahead.
- **target_update_freq** (*int*) – the target network update frequency (0 if you do not use the target network).

See also:

Please refer to *BasePolicy* for more detailed explanation.

eval () → None

Set the module in evaluation mode, except for the target network.

forward (batch: *tianshou.data.batch.Batch*, state: *Union[dict, tianshou.data.batch.Batch, numpy.ndarray, None]* = None, model: *str* = 'model', input: *str* = 'obs', eps: *Optional[float]* = None, **kwargs) → *tianshou.data.batch.Batch*

Compute action over the given batch data.

Parameters **eps** (*float*) – in [0, 1], for epsilon-greedy exploration method.

Returns

A *Batch* which has 3 keys:

- **act** the action.
- **logits** the network's raw output.
- **state** the hidden state.

See also:

Please refer to *forward()* for more detailed explanation.

learn (batch: *tianshou.data.batch.Batch*, **kwargs) → *Dict[str, float]*

Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

process_fn (batch: *tianshou.data.batch.Batch*, buffer: *tianshou.data.buffer.ReplayBuffer*, indice: *numpy.ndarray*) → *tianshou.data.batch.Batch*

Compute the n-step return for Q-learning targets:

$$G_t = \sum_{i=t}^{t+n-1} \gamma^{i-t} (1 - d_i) r_i + \gamma^n (1 - d_{t+n}) \max_a Q_{old}(s_{t+n}, \arg \max_a (Q_{new}(s_{t+n}, a)))$$

, where γ is the discount factor, $\gamma \in [0, 1]$, d_t is the done flag of step t . If there is no target network, the Q_{old} is equal to Q_{new} .

set_eps (*eps: float*) → None
Set the eps for epsilon-greedy exploration.

sync_weight () → None
Synchronize the weight for the target network.

train () → None
Set the module in training mode, except for the target network.

```
class tianshou.policy.PGPolicy(model: torch.nn.modules.module.Module, optim: torch.optim.optimizer.Optimizer, dist_fn: torch.distributions.distribution.Distribution = <class 'torch.distributions.categorical.Categorical'>, discount_factor: float = 0.99, reward_normalization: bool = False, **kwargs)
Bases: tianshou.policy.base.BasePolicy
```

Implementation of Vanilla Policy Gradient.

Parameters

- **model** (*torch.nn.Module*) – a model following the rules in *BasePolicy*. (s → logits)
- **optim** (*torch.optim.Optimizer*) – a torch.optim for optimizing the model.
- **dist_fn** (*torch.distributions.Distribution*) – for computing the action.
- **discount_factor** (*float*) – in [0, 1].

See also:

Please refer to *BasePolicy* for more detailed explanation.

forward (*batch: tianshou.data.batch.Batch, state: Union[dict, tianshou.data.batch.Batch, numpy.ndarray, None] = None, **kwargs*) → *tianshou.data.batch.Batch*
Compute action over the given batch data.

Returns

A *Batch* which has 4 keys:

- **act** the action.
- **logits** the network's raw output.
- **dist** the action distribution.
- **state** the hidden state.

See also:

Please refer to *forward()* for more detailed explanation.

learn (*batch: tianshou.data.batch.Batch, batch_size: int, repeat: int, **kwargs*) → *Dict[str, List[float]]*
Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

process_fn (*batch: tianshou.data.batch.Batch, buffer: tianshou.data.buffer.ReplayBuffer, indice: numpy.ndarray*) → *tianshou.data.batch.Batch*
Compute the discounted returns for each frame:

$$G_t = \sum_{i=t}^T \gamma^{i-t} r_i$$

, where T is the terminal time step, γ is the discount factor, $\gamma \in [0, 1]$.

```
class tianshou.policy.A2CPolicy(actor: torch.nn.modules.module.Module,
                                critic: torch.nn.modules.module.Module, op-
                                tim: torch.optim.optimizer.Optimizer, dist_fn:
                                torch.distributions.distribution.Distribution = <class
                                'torch.distributions.categorical.Categorical'>, discount_factor:
                                float = 0.99, vf_coef: float = 0.5, ent_coef: float = 0.01,
                                max_grad_norm: Optional[float] = None, gae_lambda: float
                                = 0.95, reward_normalization: bool = False, **kwargs)
```

Bases: tianshou.policy.modelfree.pg.PGPolicy

Implementation of Synchronous Advantage Actor-Critic. arXiv:1602.01783

Parameters

- **actor** (*torch.nn.Module*) – the actor network following the rules in *BasePolicy*. (s -> logits)
- **critic** (*torch.nn.Module*) – the critic network. (s -> V(s))
- **optim** (*torch.optim.Optimizer*) – the optimizer for actor and critic network.
- **dist_fn** (*torch.distributions.Distribution*) – for computing the action, defaults to *torch.distributions.Categorical*.
- **discount_factor** (*float*) – in [0, 1], defaults to 0.99.
- **vf_coef** (*float*) – weight for value loss, defaults to 0.5.
- **ent_coef** (*float*) – weight for entropy loss, defaults to 0.01.
- **max_grad_norm** (*float*) – clipping gradients in back propagation, defaults to None.
- **gae_lambda** (*float*) – in [0, 1], param for Generalized Advantage Estimation, defaults to 0.95.

See also:

Please refer to *BasePolicy* for more detailed explanation.

```
forward (batch: tianshou.data.batch.Batch, state: Union[dict, tianshou.data.batch.Batch,
numpy.ndarray, None] = None, **kwargs) → tianshou.data.batch.Batch
Compute action over the given batch data.
```

Returns

A *Batch* which has 4 keys:

- **act** the action.
- **logits** the network's raw output.
- **dist** the action distribution.
- **state** the hidden state.

See also:

Please refer to *forward()* for more detailed explanation.

```
learn (batch: tianshou.data.batch.Batch, batch_size: int, repeat: int, **kwargs) → Dict[str, List[float]]
Update policy with a given batch of data.
```

Returns A dict which includes loss and its corresponding label.

process_fn (*batch*: *tianshou.data.batch.Batch*, *buffer*: *tianshou.data.buffer.ReplayBuffer*, *indice*: *numpy.ndarray*) → *tianshou.data.batch.Batch*
 Compute the discounted returns for each frame:

$$G_t = \sum_{i=t}^T \gamma^{i-t} r_i$$

, where T is the terminal time step, γ is the discount factor, $\gamma \in [0, 1]$.

```
class tianshou.policy.DDPGPolicy (actor: torch.nn.modules.module.Module, ac-  

tor_optim: torch.optim.optimizer.Optimizer, critic:  

torch.nn.modules.module.Module, critic_optim:  

torch.optim.optimizer.Optimizer, tau: float = 0.005, gamma:  

float = 0.99, exploration_noise: float = 0.1, action_range:  

Optional[Tuple[float, float]] = None, reward_normalization:  

bool = False, ignore_done: bool = False, **kwargs)  

Bases: tianshou.policy.base.BasePolicy
```

Implementation of Deep Deterministic Policy Gradient. arXiv:1509.02971

Parameters

- **actor** (*torch.nn.Module*) – the actor network following the rules in *BasePolicy*. (s → logits)
- **actor_optim** (*torch.optim.Optimizer*) – the optimizer for actor network.
- **critic** (*torch.nn.Module*) – the critic network. (s, a → Q(s, a))
- **critic_optim** (*torch.optim.Optimizer*) – the optimizer for critic network.
- **tau** (*float*) – param for soft update of the target network, defaults to 0.005.
- **gamma** (*float*) – discount factor, in [0, 1], defaults to 0.99.
- **exploration_noise** (*float*) – the noise intensity, add to the action, defaults to 0.1.
- **action_range** ((*float*, *float*)) – the action range (minimum, maximum).
- **reward_normalization** (*bool*) – normalize the reward to Normal(0, 1), defaults to False.
- **ignore_done** (*bool*) – ignore the done flag while training the policy, defaults to False.

See also:

Please refer to *BasePolicy* for more detailed explanation.

eval () → None

Set the module in evaluation mode, except for the target network.

forward (*batch*: *tianshou.data.batch.Batch*, *state*: Union[dict, *tianshou.data.batch.Batch*, *numpy.ndarray*, None] = None, *model*: str = 'actor', *input*: str = 'obs', *eps*: Optional[float] = None, ****kwargs**) → *tianshou.data.batch.Batch*
 Compute action over the given batch data.

Parameters **eps** (*float*) – in [0, 1], for exploration use.

Returns

A *Batch* which has 2 keys:

- **act** the action.
- **state** the hidden state.

See also:

Please refer to `forward()` for more detailed explanation.

learn (*batch*: `tianshou.data.batch.Batch`, ***kwargs*) \rightarrow Dict[str, float]

Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

process_fn (*batch*: `tianshou.data.batch.Batch`, *buffer*: `tianshou.data.buffer.ReplayBuffer`, *indice*: `numpy.ndarray`) \rightarrow `tianshou.data.batch.Batch`

Pre-process the data from the provided replay buffer. Check out [Policy](#) for more information.

set_eps (*eps*: float) \rightarrow None

Set the eps for exploration.

sync_weight () \rightarrow None

Soft-update the weight for the target network.

train () \rightarrow None

Set the module in training mode, except for the target network.

```
class tianshou.policy.PPOPoly (actor: torch.nn.modules.module.Module,
                                critic: torch.nn.modules.module.Module, op-
                                tim: torch.optim.optimizer.Optimizer, dist_fn:
                                torch.distributions.distribution.Distribution, discount_factor:
                                float = 0.99, max_grad_norm: Optional[float] = None,
                                eps_clip: float = 0.2, vf_coef: float = 0.5, ent_coef: float
                                = 0.01, action_range: Optional[Tuple[float, float]] = None,
                                gae_lambda: float = 0.95, dual_clip: Optional[float] = None,
                                value_clip: bool = True, reward_normalization: bool = True,
                                **kwargs)
```

Bases: `tianshou.policy.modelfree.pg.PGPoly`

Implementation of Proximal Policy Optimization. arXiv:1707.06347

Parameters

- **actor** (`torch.nn.Module`) – the actor network following the rules in [BasePolicy](#). (s \rightarrow logits)
- **critic** (`torch.nn.Module`) – the critic network. (s \rightarrow V(s))
- **optim** (`torch.optim.Optimizer`) – the optimizer for actor and critic network.
- **dist_fn** (`torch.distributions.Distribution`) – for computing the action.
- **discount_factor** (`float`) – in [0, 1], defaults to 0.99.
- **max_grad_norm** (`float`) – clipping gradients in back propagation, defaults to None.
- **eps_clip** (`float`) – ϵ in L_{CLIP} in the original paper, defaults to 0.2.
- **vf_coef** (`float`) – weight for value loss, defaults to 0.5.
- **ent_coef** (`float`) – weight for entropy loss, defaults to 0.01.
- **action_range** ((`float`, `float`)) – the action range (minimum, maximum).
- **gae_lambda** (`float`) – in [0, 1], param for Generalized Advantage Estimation, defaults to 0.95.
- **dual_clip** (`float`) – a parameter c mentioned in arXiv:1912.09729 Equ. 5, where $c > 1$ is a constant indicating the lower bound, defaults to 5.0 (set None if you do not want to use it).

- **value_clip** (*bool*) – a parameter mentioned in arXiv:1811.02553 Sec. 4.1, defaults to `True`.
- **reward_normalization** (*bool*) – normalize the returns to `Normal(0, 1)`, defaults to `True`.

See also:

Please refer to `BasePolicy` for more detailed explanation.

forward (*batch*: `tianshou.data.batch.Batch`, *state*: `Union[dict, tianshou.data.batch.Batch, numpy.ndarray, None] = None`, ***kwargs*) → `tianshou.data.batch.Batch`
 Compute action over the given batch data.

Returns

A `Batch` which has 4 keys:

- `act` the action.
- `logits` the network's raw output.
- `dist` the action distribution.
- `state` the hidden state.

See also:

Please refer to `forward()` for more detailed explanation.

learn (*batch*: `tianshou.data.batch.Batch`, *batch_size*: *int*, *repeat*: *int*, ***kwargs*) → `Dict[str, List[float]]`
 Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

process_fn (*batch*: `tianshou.data.batch.Batch`, *buffer*: `tianshou.data.buffer.ReplayBuffer`, *indice*: `numpy.ndarray`) → `tianshou.data.batch.Batch`
 Compute the discounted returns for each frame:

$$G_t = \sum_{i=t}^T \gamma^{i-t} r_i$$

, where T is the terminal time step, γ is the discount factor, $\gamma \in [0, 1]$.

```
class tianshou.policy.TD3Policy (actor: torch.nn.modules.module.Module, ac-  
                                tor_optim: torch.optim.optimizer.Optimizer,  
                                critic1: torch.nn.modules.module.Module,  
                                critic1_optim: torch.optim.optimizer.Optimizer, critic2:  
                                torch.nn.modules.module.Module, critic2_optim:  
                                torch.optim.optimizer.Optimizer, tau: float = 0.005, gamma:  
                                float = 0.99, exploration_noise: float = 0.1, policy_noise:  
                                float = 0.2, update_actor_freq: int = 2, noise_clip: float  
                                = 0.5, action_range: Optional[Tuple[float, float]] = None,  
                                reward_normalization: bool = False, ignore_done: bool =  
                                False, **kwargs)
```

Bases: `tianshou.policy.modelfree.ddpg.DDPGPoly`

Implementation of Twin Delayed Deep Deterministic Policy Gradient, arXiv:1802.09477

Parameters

- **actor** (`torch.nn.Module`) – the actor network following the rules in `BasePolicy`. (s -> logits)
- **actor_optim** (`torch.optim.Optimizer`) – the optimizer for actor network.

- **critic1** (*torch.nn.Module*) – the first critic network. (s, a → Q(s, a))
- **critic1_optim** (*torch.optim.Optimizer*) – the optimizer for the first critic network.
- **critic2** (*torch.nn.Module*) – the second critic network. (s, a → Q(s, a))
- **critic2_optim** (*torch.optim.Optimizer*) – the optimizer for the second critic network.
- **tau** (*float*) – param for soft update of the target network, defaults to 0.005.
- **gamma** (*float*) – discount factor, in [0, 1], defaults to 0.99.
- **exploration_noise** (*float*) – the noise intensity, add to the action, defaults to 0.1.
- **policy_noise** (*float*) – the noise used in updating policy network, default to 0.2.
- **update_actor_freq** (*int*) – the update frequency of actor network, default to 2.
- **noise_clip** (*float*) – the clipping range used in updating policy network, default to 0.5.
- **action_range** ((*float*, *float*)) – the action range (minimum, maximum).
- **reward_normalization** (*bool*) – normalize the reward to Normal(0, 1), defaults to False.
- **ignore_done** (*bool*) – ignore the done flag while training the policy, defaults to False.

See also:

Please refer to [BasePolicy](#) for more detailed explanation.

eval () → None

Set the module in evaluation mode, except for the target network.

learn (*batch: tianshou.data.batch.Batch*, ***kwargs*) → Dict[str, float]

Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

sync_weight () → None

Soft-update the weight for the target network.

train () → None

Set the module in training mode, except for the target network.

```
class tianshou.policy.SACPolicy(actor: torch.nn.modules.module.Module, ac-
    tor_optim: torch.optim.optimizer.Optimizer,
    critic1: torch.nn.modules.module.Module,
    critic1_optim: torch.optim.optimizer.Optimizer, critic2:
    torch.nn.modules.module.Module, critic2_optim:
    torch.optim.optimizer.Optimizer, tau: float = 0.005,
    gamma: float = 0.99, alpha: float = 0.2, action_range:
    Optional[Tuple[float, float]] = None, reward_normalization:
    bool = False, ignore_done: bool = False, **kwargs)
```

Bases: `tianshou.policy.modelfree.ddpg.DDPGPolicy`

Implementation of Soft Actor-Critic. arXiv:1812.05905

Parameters

- **actor** (*torch.nn.Module*) – the actor network following the rules in [BasePolicy](#). (s → logits)

- **actor_optim** (*torch.optim.Optimizer*) – the optimizer for actor network.
- **critic1** (*torch.nn.Module*) – the first critic network. (s, a → Q(s, a))
- **critic1_optim** (*torch.optim.Optimizer*) – the optimizer for the first critic network.
- **critic2** (*torch.nn.Module*) – the second critic network. (s, a → Q(s, a))
- **critic2_optim** (*torch.optim.Optimizer*) – the optimizer for the second critic network.
- **tau** (*float*) – param for soft update of the target network, defaults to 0.005.
- **gamma** (*float*) – discount factor, in [0, 1], defaults to 0.99.
- **exploration_noise** (*float*) – the noise intensity, add to the action, defaults to 0.1.
- **alpha** (*float*) – entropy regularization coefficient, default to 0.2.
- **action_range** ((*float*, *float*)) – the action range (minimum, maximum).
- **reward_normalization** (*bool*) – normalize the reward to Normal(0, 1), defaults to False.
- **ignore_done** (*bool*) – ignore the done flag while training the policy, defaults to False.

See also:

Please refer to [BasePolicy](#) for more detailed explanation.

eval () → None

Set the module in evaluation mode, except for the target network.

forward (batch: *tianshou.data.batch.Batch*, state: *Union[dict, tianshou.data.batch.Batch, numpy.ndarray, None]* = None, input: *str* = 'obs', **kwargs) → *tianshou.data.batch.Batch*
 Compute action over the given batch data.

Parameters **eps** (*float*) – in [0, 1], for exploration use.

Returns

A *Batch* which has 2 keys:

- **act** the action.
- **state** the hidden state.

See also:

Please refer to [forward\(\)](#) for more detailed explanation.

learn (batch: *tianshou.data.batch.Batch*, **kwargs) → *Dict[str, float]*

Update policy with a given batch of data.

Returns A dict which includes loss and its corresponding label.

sync_weight () → None

Soft-update the weight for the target network.

train () → None

Set the module in training mode, except for the target network.

1.8 tianshou.trainer

`tianshou.trainer.gather_info` (*start_time: float, train_c: tianshou.data.collector.Collector, test_c: tianshou.data.collector.Collector, best_reward: float*) \rightarrow Dict[str, Union[float, str]]

A simple wrapper of gathering information from collectors.

Returns

A dictionary with the following keys:

- `train_step` the total collected step of training collector;
- `train_episode` the total collected episode of training collector;
- `train_time/collector` the time for collecting frames in the training collector;
- `train_time/model` the time for training models;
- `train_speed` the speed of training (frames per second);
- `test_step` the total collected step of test collector;
- `test_episode` the total collected episode of test collector;
- `test_time` the time for testing;
- `test_speed` the speed of testing (frames per second);
- `best_reward` the best reward over the test results;
- `duration` the total elapsed time.

`tianshou.trainer.test_episode` (*policy: tianshou.policy.base.BasePolicy, collector: tianshou.data.collector.Collector, test_fn: Callable[[int], None], epoch: int, n_episode: Union[int, List[int]]*) \rightarrow Dict[str, float]

A simple wrapper of testing policy in collector.

`tianshou.trainer.onpolicy_trainer` (*policy: tianshou.policy.base.BasePolicy, train_collector: tianshou.data.collector.Collector, test_collector: tianshou.data.collector.Collector, max_epoch: int, step_per_epoch: int, collect_per_step: int, repeat_per_collect: int, episode_per_test: Union[int, List[int]], batch_size: int, train_fn: Optional[Callable[[int], None]] = None, test_fn: Optional[Callable[[int], None]] = None, stop_fn: Optional[Callable[[float], bool]] = None, save_fn: Optional[Callable[[tianshou.policy.base.BasePolicy], None]] = None, log_fn: Optional[Callable[[dict], None]] = None, writer: Optional[torch.utils.tensorboard.writer.SummaryWriter] = None, log_interval: int = 1, verbose: bool = True, **kwargs*) \rightarrow Dict[str, Union[float, str]]

A wrapper for on-policy trainer procedure.

Parameters

- **policy** – an instance of the `BasePolicy` class.
- **train_collector** (`Collector`) – the collector used for training.
- **test_collector** (`Collector`) – the collector used for testing.

- **max_epoch** (*int*) – the maximum of epochs for training. The training process might be finished before reaching the `max_epoch`.
- **step_per_epoch** (*int*) – the number of step for updating policy network in one epoch.
- **collect_per_step** (*int*) – the number of frames the collector would collect before the network update. In other words, collect some frames and do one policy network update.
- **repeat_per_collect** (*int*) – the number of repeat time for policy learning, for example, set it to 2 means the policy needs to learn each given batch data twice.
- **episode_per_test** (*int or list of ints*) – the number of episodes for one policy evaluation.
- **batch_size** (*int*) – the batch size of sample data, which is going to feed in the policy network.
- **train_fn** (*function*) – a function receives the current number of epoch index and performs some operations at the beginning of training in this epoch.
- **test_fn** (*function*) – a function receives the current number of epoch index and performs some operations at the beginning of testing in this epoch.
- **save_fn** (*function*) – a function for saving policy when the undiscounted average mean reward in evaluation phase gets better.
- **stop_fn** (*function*) – a function receives the average undiscounted returns of the testing result, return a boolean which indicates whether reaching the goal.
- **log_fn** (*function*) – a function receives env info for logging.
- **writer** (*torch.utils.tensorboard.SummaryWriter*) – a TensorBoard SummaryWriter.
- **log_interval** (*int*) – the log interval of the writer.
- **verbose** (*bool*) – whether to print the information.

Returns See `gather_info()`.

```
tianshou.trainer.offpolicy_trainer(policy: tianshou.policy.base.BasePolicy, train_collector:
                                     tianshou.data.collector.Collector, test_collector:
                                     tianshou.data.collector.Collector, max_epoch:
                                     int, step_per_epoch: int, collect_per_step: int,
                                     episode_per_test: Union[int, List[int]], batch_size:
                                     int, train_fn: Optional[Callable[[int], None]] = None,
                                     test_fn: Optional[Callable[[int], None]] = None, stop_fn:
                                     Optional[Callable[[float], bool]] = None, save_fn:
                                     Optional[Callable[[tianshou.policy.base.BasePolicy],
                                     None]] = None, log_fn: Optional[Callable[[dict], None]] = None, writer: Op-
                                     tional[torch.utils.tensorboard.writer.SummaryWriter]
                                     = None, log_interval: int = 1, verbose: bool = True,
                                     **kwargs) → Dict[str, Union[float, str]]
```

A wrapper for off-policy trainer procedure.

Parameters

- **policy** – an instance of the `BasePolicy` class.
- **train_collector** (*Collector*) – the collector used for training.
- **test_collector** (*Collector*) – the collector used for testing.

- **max_epoch** (*int*) – the maximum of epochs for training. The training process might be finished before reaching the `max_epoch`.
- **step_per_epoch** (*int*) – the number of step for updating policy network in one epoch.
- **collect_per_step** (*int*) – the number of frames the collector would collect before the network update. In other words, collect some frames and do one policy network update.
- **episode_per_test** – the number of episodes for one policy evaluation.
- **batch_size** (*int*) – the batch size of sample data, which is going to feed in the policy network.
- **train_fn** (*function*) – a function receives the current number of epoch index and performs some operations at the beginning of training in this epoch.
- **test_fn** (*function*) – a function receives the current number of epoch index and performs some operations at the beginning of testing in this epoch.
- **save_fn** (*function*) – a function for saving policy when the undiscounted average mean reward in evaluation phase gets better.
- **stop_fn** (*function*) – a function receives the average undiscounted returns of the testing result, return a boolean which indicates whether reaching the goal.
- **log_fn** (*function*) – a function receives env info for logging.
- **writer** (*torch.utils.tensorboard.SummaryWriter*) – a TensorBoard SummaryWriter.
- **log_interval** (*int*) – the log interval of the writer.
- **verbose** (*bool*) – whether to print the information.

Returns See `gather_info()`.

1.9 tianshou.exploration

class `tianshou.exploration.OUNoise` (*sigma: float = 0.3, theta: float = 0.15, dt: float = 0.01, x0: Union[float, numpy.ndarray, None] = None*)

Bases: `object`

Class for Ornstein-Uhlenbeck process, as used for exploration in DDPG. Usage:

```
# init
self.noise = OUNoise()
# generate noise
noise = self.noise(logits.shape, eps)
```

For required parameters, you can refer to the stackoverflow page. However, our experiment result shows that (similar to OpenAI SpinningUp) using vanilla gaussian process has little difference from using the Ornstein-Uhlenbeck process.

__call__ (*size: tuple, mu: float = 0.1*) → `numpy.ndarray`

Generate new noise. Return a `numpy.ndarray` which size is equal to *size*.

reset () → `None`

Reset to the initial state.

1.10 tianshou.utils

class tianshou.utils.**MovAvg** (*size: int = 100*)

Bases: object

Class for moving average. It will automatically exclude the infinity and NaN. Usage:

```
>>> stat = MovAvg(size=66)
>>> stat.add(torch.tensor(5))
5.0
>>> stat.add(float('inf')) # which will not add to stat
5.0
>>> stat.add([6, 7, 8])
6.5
>>> stat.get()
6.5
>>> print(f'{stat.mean():.2f}±{stat.std():.2f}')
6.50±1.12
```

add (*x: Union[float, list, numpy.ndarray, torch.Tensor]*) → float

Add a scalar into *MovAvg*. You can add `torch.Tensor` with only one element, a python scalar, or a list of python scalar.

get () → float

Get the average.

mean () → float

Get the average. Same as *get* ().

std () → float

Get the standard deviation.

1.11 Contributing to Tianshou

1.11.1 Install Develop Version

To install Tianshou in an “editable” mode, run

```
pip3 install -e ".[dev]"
```

in the main directory. This installation is removable by

```
python3 setup.py develop --uninstall
```

1.11.2 PEP8 Code Style Check

We follow PEP8 python code style. To check, in the main directory, run:

```
flake8 . --count --show-source --statistics
```

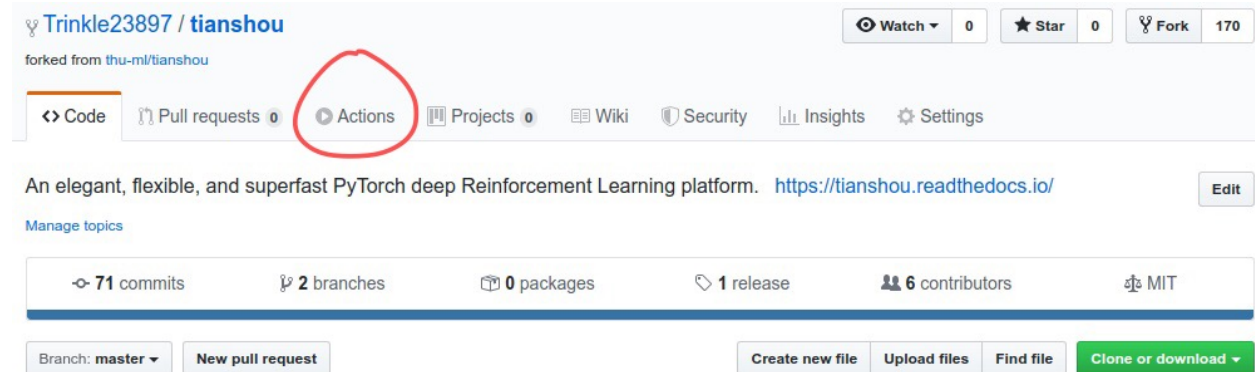
1.11.3 Test Locally

This command will run automatic tests in the main directory

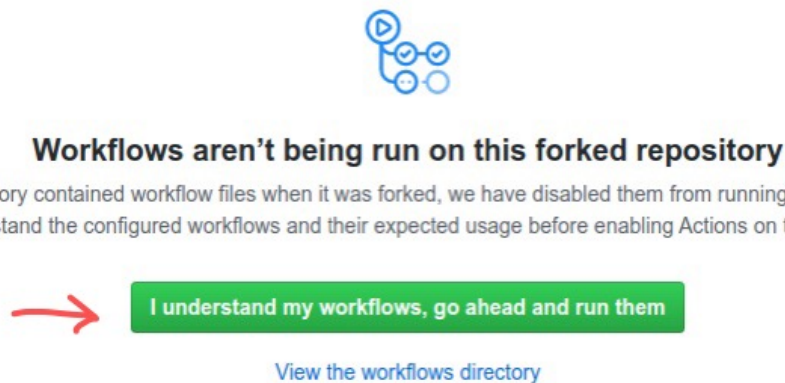
```
pytest test --cov tianshou -s --durations 0 -v
```

1.11.4 Test by GitHub Actions

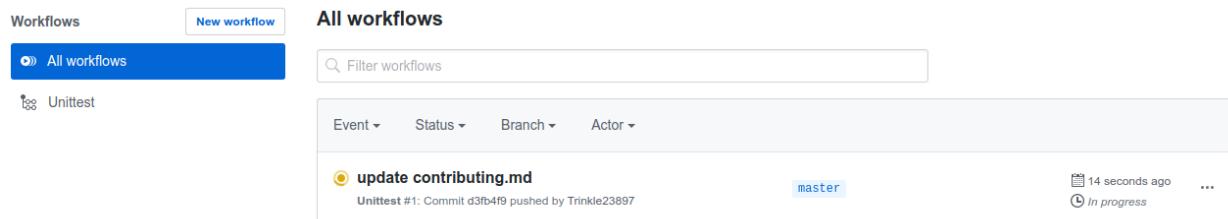
1. Click the Actions button in your own repo:



2. Click the green button:



3. You will see Actions Enabled. on the top of html page.
4. When you push a new commit to your own repo (e.g. `git push`), it will automatically run the test in this page:



1.11.5 Documentation

Documentations are written under the `docs/` directory as ReStructuredText (`.rst`) files. `index.rst` is the main page. A Tutorial on ReStructuredText can be found [here](#).

API References are automatically generated by [Sphinx](#) according to the outlines under `docs/api/` and should be modified when any code changes.

To compile documentation into webpages, run

```
make html
```

under the `docs/` directory. The generated webpages are in `docs/_build` and can be viewed with browsers.

1.11.6 Chinese Documentation

Chinese documentation is in https://github.com/thu-ml/tianshou-docs-zh_CN.

1.12 Contributor

We always welcome contributions to help make Tianshou better. Below are an incomplete list of our contributors (find more on [this page](#)).

- Jiayi Weng ([Trinkle23897](#))
- Minghao Zhang ([Mehooz](#))

INDICES AND TABLES

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

BIBLIOGRAPHY

- [MKS+15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015. URL: <https://doi.org/10.1038/nature14236>, doi:10.1038/nature14236.
- [LHP+16] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*. 2016. URL: <http://arxiv.org/abs/1509.02971>.
- [SWD+17] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, 2017. URL: <http://arxiv.org/abs/1707.06347>, arXiv:1707.06347.

PYTHON MODULE INDEX

t

- `tianshou.data`, [15](#)
- `tianshou.env`, [22](#)
- `tianshou.exploration`, [38](#)
- `tianshou.policy`, [25](#)
- `tianshou.trainer`, [36](#)
- `tianshou.utils`, [39](#)

Symbols

`__call__()` (*tianshou.exploration.OUNoise method*), 38
`__getitem__()` (*tianshou.data.Batch method*), 16
`__getitem__()` (*tianshou.data.PrioritizedReplayBuffer method*), 19
`__getitem__()` (*tianshou.data.ReplayBuffer method*), 18
`__len__()` (*tianshou.data.Batch method*), 16
`__len__()` (*tianshou.data.ReplayBuffer method*), 18
`__len__()` (*tianshou.env.BaseVectorEnv method*), 22

A

`A2CPolicy` (*class in tianshou.policy*), 29
`add()` (*tianshou.data.PrioritizedReplayBuffer method*), 19
`add()` (*tianshou.data.ReplayBuffer method*), 18
`add()` (*tianshou.utils.MovAvg method*), 39
`append()` (*tianshou.data.Batch method*), 16

B

`BasePolicy` (*class in tianshou.policy*), 25
`BaseVectorEnv` (*class in tianshou.env*), 22
`Batch` (*class in tianshou.data*), 15

C

`close()` (*tianshou.data.Collector method*), 21
`close()` (*tianshou.env.BaseVectorEnv method*), 22
`close()` (*tianshou.env.RayVectorEnv method*), 25
`close()` (*tianshou.env.SubprocVectorEnv method*), 24
`close()` (*tianshou.env.VectorEnv method*), 23
`collect()` (*tianshou.data.Collector method*), 21
`Collector` (*class in tianshou.data*), 19
`compute_episodic_return()` (*tianshou.policy.BasePolicy static method*), 26

D

`DDPGPolicy` (*class in tianshou.policy*), 31
`DQNPolicy` (*class in tianshou.policy*), 27

E

`eval()` (*tianshou.policy.DDPGPolicy method*), 31
`eval()` (*tianshou.policy.DQNPolicy method*), 28
`eval()` (*tianshou.policy.SACPolicy method*), 35
`eval()` (*tianshou.policy.TD3Policy method*), 34

F

`forward()` (*tianshou.policy.A2CPolicy method*), 30
`forward()` (*tianshou.policy.BasePolicy method*), 26
`forward()` (*tianshou.policy.DDPGPolicy method*), 31
`forward()` (*tianshou.policy.DQNPolicy method*), 28
`forward()` (*tianshou.policy.ImitationPolicy method*), 27
`forward()` (*tianshou.policy.PGPolicy method*), 29
`forward()` (*tianshou.policy.PPOPolicy method*), 33
`forward()` (*tianshou.policy.SACPolicy method*), 35

G

`gather_info()` (*in module tianshou.trainer*), 36
`get()` (*tianshou.data.Batch method*), 16
`get()` (*tianshou.data.ReplayBuffer method*), 18
`get()` (*tianshou.utils.MovAvg method*), 39
`get_env_num()` (*tianshou.data.Collector method*), 21

I

`ImitationPolicy` (*class in tianshou.policy*), 27

K

`keys()` (*tianshou.data.Batch method*), 16

L

`learn()` (*tianshou.policy.A2CPolicy method*), 30
`learn()` (*tianshou.policy.BasePolicy method*), 26
`learn()` (*tianshou.policy.DDPGPolicy method*), 32
`learn()` (*tianshou.policy.DQNPolicy method*), 28
`learn()` (*tianshou.policy.ImitationPolicy method*), 27
`learn()` (*tianshou.policy.PGPolicy method*), 29
`learn()` (*tianshou.policy.PPOPolicy method*), 33
`learn()` (*tianshou.policy.SACPolicy method*), 35
`learn()` (*tianshou.policy.TD3Policy method*), 34
`ListReplayBuffer` (*class in tianshou.data*), 19

M

`mean()` (*tianshou.utils.MovAvg method*), 39

module

`tianshou.data`, 15

`tianshou.env`, 22

`tianshou.exploration`, 38

`tianshou.policy`, 25

`tianshou.trainer`, 36

`tianshou.utils`, 39

`MovAvg` (*class in tianshou.utils*), 39

O

`offpolicy_trainer()` (*in module tianshou.trainer*), 37

`onpolicy_trainer()` (*in module tianshou.trainer*), 36

`OUNoise` (*class in tianshou.exploration*), 38

P

`PGPolicy` (*class in tianshou.policy*), 29

`PPOPolicy` (*class in tianshou.policy*), 32

`PrioritizedReplayBuffer` (*class in tianshou.data*), 19

`process_fn()` (*tianshou.policy.A2CPolicy method*), 30

`process_fn()` (*tianshou.policy.BasePolicy method*), 27

`process_fn()` (*tianshou.policy.DDPGPolicy method*), 32

`process_fn()` (*tianshou.policy.DQNPolicy method*), 28

`process_fn()` (*tianshou.policy.PGPolicy method*), 29

`process_fn()` (*tianshou.policy.PPOPolicy method*), 33

R

`RayVectorEnv` (*class in tianshou.env*), 24

`render()` (*tianshou.data.Collector method*), 21

`render()` (*tianshou.env.BaseVectorEnv method*), 22

`render()` (*tianshou.env.RayVectorEnv method*), 25

`render()` (*tianshou.env.SubprocVectorEnv method*), 24

`render()` (*tianshou.env.VectorEnv method*), 23

`ReplayBuffer` (*class in tianshou.data*), 17

`reset()` (*tianshou.data.Collector method*), 21

`reset()` (*tianshou.data.ListReplayBuffer method*), 19

`reset()` (*tianshou.data.PrioritizedReplayBuffer method*), 19

`reset()` (*tianshou.data.ReplayBuffer method*), 18

`reset()` (*tianshou.env.BaseVectorEnv method*), 22

`reset()` (*tianshou.env.RayVectorEnv method*), 25

`reset()` (*tianshou.env.SubprocVectorEnv method*), 24

`reset()` (*tianshou.env.VectorEnv method*), 23

`reset()` (*tianshou.exploration.OUNoise method*), 38

`reset_buffer()` (*tianshou.data.Collector method*), 21

`reset_env()` (*tianshou.data.Collector method*), 21

S

`SACPolicy` (*class in tianshou.policy*), 34

`sample()` (*tianshou.data.Collector method*), 22

`sample()` (*tianshou.data.PrioritizedReplayBuffer method*), 19

`sample()` (*tianshou.data.ReplayBuffer method*), 18

`seed()` (*tianshou.data.Collector method*), 22

`seed()` (*tianshou.env.BaseVectorEnv method*), 22

`seed()` (*tianshou.env.RayVectorEnv method*), 25

`seed()` (*tianshou.env.SubprocVectorEnv method*), 24

`seed()` (*tianshou.env.VectorEnv method*), 23

`set_eps()` (*tianshou.policy.DDPGPolicy method*), 32

`set_eps()` (*tianshou.policy.DQNPolicy method*), 28

`split()` (*tianshou.data.Batch method*), 16

`std()` (*tianshou.utils.MovAvg method*), 39

`step()` (*tianshou.env.BaseVectorEnv method*), 22

`step()` (*tianshou.env.RayVectorEnv method*), 25

`step()` (*tianshou.env.SubprocVectorEnv method*), 24

`step()` (*tianshou.env.VectorEnv method*), 23

`SubprocVectorEnv` (*class in tianshou.env*), 24

`sync_weight()` (*tianshou.policy.DDPGPolicy method*), 32

`sync_weight()` (*tianshou.policy.DQNPolicy method*), 29

`sync_weight()` (*tianshou.policy.SACPolicy method*), 35

`sync_weight()` (*tianshou.policy.TD3Policy method*), 34

T

`TD3Policy` (*class in tianshou.policy*), 33

`test_episode()` (*in module tianshou.trainer*), 36

`tianshou.data`

module, 15

`tianshou.env`

module, 22

`tianshou.exploration`

module, 38

`tianshou.policy`

module, 25

`tianshou.trainer`

module, 36

`tianshou.utils`

module, 39

`to_numpy()` (*in module tianshou.data*), 16

`to_numpy()` (*tianshou.data.Batch method*), 16

`to_torch()` (*in module tianshou.data*), 16

`to_torch()` (*tianshou.data.Batch method*), 16

`train()` (*tianshou.policy.DDPGPolicy method*), 32

`train()` (*tianshou.policy.DQNPolicy method*), 29

`train()` (*tianshou.policy.SACPolicy method*), [35](#)
`train()` (*tianshou.policy.TD3Policy method*), [34](#)

U

`update()` (*tianshou.data.ReplayBuffer method*), [18](#)
`update_weight()` (*tianshou.data.PrioritizedReplayBuffer method*),
[19](#)

V

`values()` (*tianshou.data.Batch method*), [16](#)
`VectorEnv` (*class in tianshou.env*), [23](#)