

EFFICIENT REINFORCEMENT LEARNING EXPERIMENTATION IN PYTORCH

Anonymous authors

Paper under double-blind review

1 SUPPLEMENTARY MATERIAL

The online documentation of the package can be found at: To be disclosed

1.1 EXPERIMENTS ON ATARI 2600 ENVIRONMENTS

We used the ensuing parameters, following Schulman et al. (2017). We also used the same CNN network architecture:

Table 1: PPO + GAE Agent - hyperparameters values.

PARAMETER	VALUE
ADAM LR	2.5E-4
CLIP-PARAM	0.15
GAMMA	0.99
FRAME-SKIP	4
FRAME-STACK	4
NUM-STEPS	128
NUM-MINI-BATCH	4
NUM-EPOCHS	3
ENTROPY-COEF	0.01
VALUE-LOSS-COEF	1.0
ENVIRONMENT VECTOR LENGTH	8
GAE-LAMBDA	0.95
MAX-GRAD-NORM	0.5

The experiment was executed with the following training script. During training, the learning rate and the PPO clipping parameter are decayed also following Schulman et al. (2017):

```

1 import sys
2 import ray
3 import numpy as np
4
5 from xxxxxxxx import utils
6 from xxxxxxxx import Learner
7 from xxxxxxxx.scheme import Scheme
8 from xxxxxxxx.agent.algorithms import PPO
9 from xxxxxxxx.agent.env import VecEnv
10 from xxxxxxxx.agent.storages import GAEBuffer
11 from xxxxxxxx.agent.actors import OnPolicyActor, get_feature_extractor
12 from xxxxxxxx.envs import atari_train_env_factory
13
14 seed = 1
15 env_id = "PongNoFrameskip-v4"
16 log_dir = "/tmp/atari/{}/{}".format(env_id, seed)
17 col_workers_communication = "synchronous"
18 grad_workers_communication = "synchronous"
19
20 if __name__ == "__main__":
21
22     utils.cleanup_log_dir(log_dir)
23     ray.init()
24
25     # 1. Define Train Vector of Envs
26     train_envs_factory, action_space, obs_space = \
27         VecEnv.create_factory(
28             env_fn=atari_train_env_factory,
29             env_kwargs={
30                 "env_id": env_id,

```

```

31         "frame_stack": 4,
32         "seed": seed),
33         vec_env_size=8, log_dir=log_dir,
34         info_keywords=('rr', 'rrr', 'lives'))
35
36     # 2. Define RL training algorithm
37     algo_factory, algo_name = PPO.create_factory(
38         lr=2.5e-4, num_epochs=3, clip_param=0.1,
39         entropy_coeff=0.01, value_loss_coeff=1.0,
40         max_grad_norm=0.5, num_mini_batch=4,
41         gamma=0.99)
42
43     # 3. Define RL Policy
44     actor_factory = OnPolicyActor.create_factory(
45         obs_space, action_space, algo_name,
46         get_feature_extractor("CNN"))
47
48     # 4. Define rollouts storage
49     storage_factory = GAEBuffer.create_factory(
50         size=128, gae_lambda=0.95)
51
52     # 5. Define scheme
53     params = {}
54
55     # add core modules
56     params.update({
57         "algo_factory": algo_factory,
58         "actor_factory": actor_factory,
59         "storage_factory": storage_factory,
60         "train_envs_factory": train_envs_factory})
61
62     # add collection specs
63     params.update({
64         "num_col_workers": 1,
65         "col_workers_communication": col_workers_communication,
66         "col_workers_resources": {
67             "num_cpus": 8, "num_gpus": 0.5}})
68
69     # add gradient specs
70     params.update({
71         "num_grad_workers": 1,
72         "grad_workers_communication": grad_workers_communication,
73         "grad_workers_resources": {
74             "num_cpus": 1.0, "num_gpus": 0.5}})
75
76     scheme = Scheme(**params)
77
78     # 6. Define learner
79     learner = Learner(
80         scheme, target_steps=10000000, log_dir=log_dir)
81
82     # 7. Define train loop
83     updates = 0
84     iterations = 0
85
86     # lr and clip parameters decay values
87     total_updates = np.ceil((10000000 * 3 * 4) / (8 * 128))
88     alpha = np.linspace(1.0, 0.0, int(total_updates))
89
90     while not learner.done():
91
92         learner.step()
93         learner.print_info()
94
95         if iterations % 100 == 0:
96             save_name = learner.save_model()
97
98         iterations += 1
99         learner.update_algo_parameter(
100             "lr", alpha[updates] * 2.5e-4)
101         learner.update_algo_parameter(
102             "clip_param", alpha[updates] * 0.1)
103
104

```

Listing 1: Training script Atari 2600

1.2 EXPERIMENTS ON PYBULLET ENVIRONMENTS

We used the following parameters, hand-picked to perform well on the single-threaded regime. As a feature extractor, we used a multilayer perceptron with two hidden layers of 256 each neurons and ReLU activation function:

Table 2: SAC Agent - hyperparameters values.

PARAMETER	VALUE
ADAM LR Q NETWORKS	1E-3
ADAM LR POLICY	1E-4
ADAM LR ALPHA	1E-05
GAMMA	0.98
INITIAL ALPHA	0.2
ENVIRONMENT VECTOR LENGTH	1
POLIAK UPDATE	0.995
BUFFER SIZE	1500000
MINI-BATCH SIZE	256
START STEPS	10000
NUM UPDATES	32
UPDATEEVERY	128

We used the following training script:

```

1 import ray
2 import time
3 from xxxxxxxxx import utils
4 from xxxxxxxxx import Learner
5 from xxxxxxxxx.scheme import Scheme
6 from xxxxxxxxx.agent.algos import SAC
7 from xxxxxxxxx.agent.env import VecEnv
8 from xxxxxxxxx.agent.storages import ReplayBuffer
9 from xxxxxxxxx.envs import pybullet_train_env_factory
10 from xxxxxxxxx.agent.actors import OffPolicyActor, get_feature_extractor
11
12 seed = 1
13 max_time = 3600 * 8 # 8h
14 env_id = "AntBulletEnv-v0"
15 log_dir = "/tmp/pybullet/{}({})".format(env_id, seed)
16 num_col_workers = 1
17 num_grad_workers = 1
18 col_workers_communication = "synchronous"
19 grad_workers_communication = "synchronous"
20
21 if __name__ == "__main__":
22
23     utils.cleanup_log_dir(log_dir)
24     ray.init()
25
26     # Define Environment Vector
27     train_envs_factory, action_space, obs_space = \
28         VecEnv.create_factory(
29             vec_env_size=1,
30             log_dir=log_dir,
31             env_fn=pybullet_train_env_factory,
32             env_kwargs={"env_id": env_id, "seed": seed})
33
34     # Define RL training algorithm
35     algo_factory, algo_name = SAC.create_factory(
36         lr_pi=0.0001, lr_q=0.001, lr_alpha=1e-05,
37         initial_alpha=0.1, gamma=0.98,
38         polyak=0.995, num_updates=8,
39         update_every=128, start_steps=5000,
40         mini_batch_size=400)
41
42     # Define RL Actor
43     actor_factory = OffPolicyActor.create_factory(
44         obs_space, action_space, algo_name,
45         get_feature_extractor("MLP"))
46
47     # Define rollouts storage
48     storage_factory = ReplayBuffer.create_factory(
49         size=15000000)
50
51     # Core components params
52     scheme_parameters = {
53         "algo_factory": algo_factory,
54         "actor_factory": actor_factory,
55         "storage_factory": storage_factory,
56         "train_envs_factory": train_envs_factory}
57
58     # Collection operation params
59     scheme_parameters.update({
60         "num_col_workers": num_col_workers,
61         "col_workers_communication":
62             col_workers_communication})
63

```

```

64     # Gradient computation operation params
65     scheme_parameters.update({
66         "num_grad_workers": num_grad_workers,
67         "grad_workers_communication": grad_workers_communication})
68
69     scheme = Scheme(**scheme_parameters)
70
71
72     # 6. Define learner
73     learner = Learner(
74         scheme, target_steps=10000000, log_dir=log_dir)
75
76     # 7. Define train loop
77     iterations = 0
78     start_time = time.time()
79     while not learner.done():
80         learner.step()
81         if iterations % 1 == 0:
82             learner.print_info()
83         if iterations % 250000 == 0:
84             learner.save_model()
85         if max_time != -1 and (
86             time.time() - start_time) > max_time:
87             break
88         iterations += 1

```

Listing 2: Training script PyBullet

1.3 EXPERIMENTS ON OBSTACLE TOWER UNITY3D CHALLENGE ENVIRONMENT

Table 3: Classification accuracies for naive Bayes and flexible Bayes on various data sets.

PARAMETER	VALUE
ADAM LR	4E-4
CLIP-PARAM	0.15
GAMMA	0.99
FRAME-SKIP	2
FRAME-STACK	4
NUM-STEPS	800
NUM-MINI-BATCH	6
NUM-EPOCHS	4
ENTROPY-COEF	0.01
VALUE-LOSS-COEF	0.2
NUM-WORKERS	1
ENVIRONMENT VECTOR LENGTH	16
GAE-LAMBDA	0.95
MAX-GRAD-NORM	0.5
RECURRENT NETWORK	TRUE
RECURRENT HIDDEN SIZE	256

We use the network architecture proposed in Espeholt et al. (2018) but we initialize its weights according to Fixup Zhang et al. (2019). We end our network with a gated recurrent unit (GRU) Cho et al. (2014) with a hidden layer of size 256 neurons. Gradients are computed using Adam optimizer Kingma & Ba (2014), with a starting learning rate of 4e-4 decayed by a factor of 0.25 both after 100 million steps and 400 million steps. The value coefficient, the entropy coefficient and the clipping parameters of the PPO algorithm are set to 0.2, 0.01 and 0.15 respectively. We use a discount factor gamma of 0.99. Furthermore, the horizon parameters is set to 800 steps and rollout collections are parallelized using environment vectors of size 16. Gradients are computed in minibatches of size 1600 for 2 epochs. Finally we also use generalized advantage estimation Schulman et al. (2015) with lambda 0.95. We use frame skip 2 and frame stack 4. We restart each new episode at a randomly selected floor between 0 and the higher floor reached in the previous episode.

We reshape the reward received from the environment in order to increase training efficiency. The reward upon the agent completing a floor is +1, and +0.1 is provided for opening doors, solving puzzles, or picking up keys. We additionally reward the agent with an extra +1 to pick up keys, +0.002 to detect boxes, +0.001 to find box intended locations, +1.5 to place the boxes target locations and +0.002 for any event that increases remaining time. We also reduce the action set from

the initial 54 actions to 6 (rotate camera in both directions, go forward, and go forward while turning left, turning right or jumping).

We used the following training script:

```

1 import time
2 from xxxxxxxxx import Learner
3 from xxxxxxxxx.scheme import Scheme
4 from xxxxxxxxx.agent.algorithms import PPO
5 from xxxxxxxxx.agent.env import VecEnv
6 from xxxxxxxxx.agent.storages import GAEBuffer
7 from xxxxxxxxx.envs import obstacle_train_env_factory
8 from xxxxxxxxx.agent.actors import OnPolicyActor, \
9     get_feature_extractor
10
11 # Define Train Vector of Envs
12 envs_factory, action_space, obs_space = \
13     VecEnv.create_factory(
14         env_fn=obstacle_train_env_factory,
15         env_kwargs={"frame_skip": 2, "frame_stack": 4},
16         log_dir='/tmp/obstacle_tower_agent',
17         vec_env_size=8, info_keywords=('floor', 'start', 'seed'))
18
19 # Define RL training algorithm
20 algo_factory, algo_name = PPO.create_factory(
21     lr=2.5e-5, num_epochs=2, clip_param=0.15,
22     entropy_coef=0.01, value_loss_coef=0.2,
23     max_grad_norm=0.5, num_mini_batch=8,
24     use_clipped_value_loss=True, gamma=0.99)
25
26 # Define RL Policy
27 actor_factory = OnPolicyActor.create_factory(
28     obs_space, action_space, algo_name,
29     get_feature_extractor("fixup"),
30     recurrent_policy=True,
31     recurrent_hidden_size=256)
32
33 # Define rollouts storage
34 storage_factory = GAEBuffer.create_factory(
35     size=800, gae_lambda=0.95)
36
37 # 6. Define scheme
38 params = {}
39
40 # add core modules
41 params.update({
42     "algo_factory": algo_factory,
43     "actor_factory": actor_factory,
44     "storage_factory": storage_factory,
45     "train_envs_factory": envs_factory})
46
47 # add collection specs
48 params.update({
49     "num_col_workers": 4,
50     "col_communication": "asynchronous",
51     "col_worker_resources": {
52         "num_cpus": 1, "num_gpus": 0.5,
53         "object_store_memory": 2.5 * 1024 ** 3,
54         "memory": 2.5 * 1024 ** 3}})
55
56 # add gradient specs
57 # num_cpus is equal to 32 - 4 to force ray to locate
58 # 4 collection workers and 1 gradient worker in each
59 # machine with 32 cpus
60 params.update({
61     "num_grad_workers": 2,
62     "grad_communication": "synchronous",
63     "grad_worker_resources": {
64         "num_cpus": 32 - 4,
65         "num_gpus": 1.0,
66         "object_store_memory": 2.5 * 1024 ** 3,
67         "memory": 2.5 * 1024 ** 3}})
68
69 scheme = Scheme(**params)
70
71 # 7. Define learner
72 learner = Learner(
73     scheme, target_steps=600000000,
74     log_dir='/tmp/obstacle_tower_agent')
75
76 # 8. Define train loop
77 iterations = 0
78 num_lr_updates = 0
79 start_time = time.time()
80 while not learner.done():
81
82     learner.step()
83     learner.print_info()
84
85     if iterations % 100 == 0:
86         save_name = learner.save_model()

```

```

87
88     if learner.num_samples_collected > 100000000 and num_lr_updates == 0:
89         learner.update_algo_parameter("lr", 4e-4 * 0.25)
90         num_lr_updates += 1
91
92     elif learner.num_samples_collected > 400000000 and num_lr_updates == 1:
93         learner.update_algo_parameter("lr", 4e-4 * 0.25 ** 2)
94         num_lr_updates += 1
95
96     iterations += 1

```

Listing 3: Training script Obstacle Tower Unity3D Challenge

1.4 SCHEME CONFIGURATION OF POPULAR RL AGENTS

Table 4 shows the architecture configurations used by some popular RL Agents with which we have experimented in this work. All these architectures can be configured via our *Scheme* class.

Table 4: *Scheme* configurations to match the underlying training architectures of some popular RL Agents.

AGENT ARCHITECTURE	NUMBER OF GRADIENT WORKERS	NUMBER OF COLLECTION WORKER (PER EACH GRADIENT WORKER)	GRADIENT WORKERS COMMUNICATION	COLLECTION WORKERS COMMUNICATION
SINGLE THREADED	1	1	SYNCHRONOUS	SYNCHRONOUS
DPPO	≥ 1	1	SYNCHRONOUS	SYNCHRONOUS
APPO	≥ 1	1	ASYNCHRONOUS	SYNCHRONOUS
IMPALA / APEX	1	≥ 1	SYNCHRONOUS	ASYNCHRONOUS
RAPID	≥ 1	≥ 1	SYNCHRONOUS	ASYNCHRONOUS
ASYNC RAPID	≥ 1	≥ 1	ASYNCHRONOUS	ASYNCHRONOUS

REFERENCES

- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymir Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. *arXiv preprint arXiv:1802.01561*, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Hongyi Zhang, Yann N Dauphin, and Tengyu Ma. Fixup initialization: Residual learning without normalization. *arXiv preprint arXiv:1901.09321*, 2019.