

```

1 # Source: https://github.com/vwxyzjn/cleanrl/blob/master/cleanrl/dqn\_atari.py
2 # docs and experiment results can be found at https://docs.cleanrl.dev/rl-algorithms/dqn/#dqn\_ataripy
3
4 import os
5 import random
6 import time
7 from dataclasses import dataclass
8
9 import gymnasium as gym
10 import numpy as np
11 import torch
12 import torch.nn as nn
13 import torch.nn.functional as F
14 import torch.optim as optim
15 import tyro
16 from stable_baselines3.common.atari_wrappers import (
17     ClipRewardEnv,
18     EpisodicLifeEnv,
19     FireResetEnv,
20     MaxAndSkipEnv,
21     NoopResetEnv,
22 )
23 from stable_baselines3.common.buffer import ReplayBuffer
24 from torch.utils.tensorboard import SummaryWriter
25
26
27 #####
28 # Instantiate a config object (for tyro) #
29 #####
30
31 @dataclass
32 class Args:
33     exp_name: str = os.path.basename(__file__)[:-len(".py")]
34     """the name of this experiment"""
35     seed: int = 1
36     """seed of the experiment"""
37     torch_deterministic: bool = True
38     """if toggled, `torch.backends.cudnn.deterministic=False`"""
39     cuda: bool = True
40     """if toggled, cuda will be enabled by default"""
41     track: bool = False
42     """if toggled, this experiment will be tracked with Weights and Biases"""
43     wandb_project_name: str = "cleanRL"
44     """the wandb's project name"""
45     wandb_entity: str = None
46     """the entity (team) of wandb's project"""
47     capture_video: bool = False
48     """whether to capture videos of the agent performances (check out `videos` folder)"""
49     save_model: bool = False
50     """whether to save model into the `runs/{run_name}` folder"""
51     upload_model: bool = False
52     """whether to upload the saved model to huggingface"""
53     hf_entity: str = ""
54     """the user or org name of the model repository from the Hugging Face Hub"""
55
56     # Algorithm specific arguments
57     env_id: str = "BreakoutNoFrameskip-v4"
58     """the id of the environment"""
59     total_timesteps: int = 1000000
60     """total timesteps of the experiments"""
61     learning_rate: float = 1e-4
62     """the learning rate of the optimizer"""
63     num_envs: int = 1
64     """the number of parallel game environments"""
65     buffer_size: int = 100000
66     """the replay memory buffer size"""
67     gamma: float = 0.99
68     """the discount factor gamma"""
69     tau: float = 1.0
70     """the target network update rate"""
71     target_network_frequency: int = 1000
72     """the timesteps it takes to update the target network"""
73     batch_size: int = 32
74     """the batch size of sample from the replay memory"""
75     start_e: float = 1
76     """the starting epsilon for exploration"""
77     end_e: float = 0.01
78     """the ending epsilon for exploration"""
79     exploration_fraction: float = 0.10
80     """the fraction of `total-timesteps` it takes from start-e to go end-e"""
81     learning_starts: int = 80000

```

```

82     """timestep to start learning"""
83     train_frequency: int = 4
84     """the frequency of training"""
85
86
87 def make_env(env_id, seed, idx, capture_video, run_name):
88     """
89     Make environment
90
91     Args:
92     env_id (str): ID of the environment (refer to the environment documentations). \n
93     seed (int): Seed in generating the environment. \n
94     idx (int) : Index of the environment. \n
95     capture_video (boolean): A flag that decides whether to create an environment that will record episodes. \n
96     run_name (str): The name of the folder where the recorded videos will be saved to. \n
97
98     Returns:
99     thunk (func): A function that returns a created environments.
100    """
101    def thunk():
102        # Create an environment with render_mode
103        # if capture_video is True and only for the first environment
104        # Record video: Record video intermittently at episode intervals
105        # (https://gymnasium.farama.org/api/wrappers/misc\_wrappers/#gymnasium.wrappers.RecordVideo)
106        if capture_video and idx == 0:
107            env = gym.make(env_id, render_mode="rgb_array")
108            env = gym.wrappers.RecordVideo(env, f"videos/{run_name}")
109        else:
110            env = gym.make(env_id)
111
112        # Enable the environment to keep track of cumulative rewards and episode lengths.
113        # Save into "info" at the end of the episodes
114        # (https://gymnasium.farama.org/api/wrappers/misc\_wrappers/#gymnasium.wrappers.RecordEpisodeStatistics)
115        env = gym.wrappers.RecordEpisodeStatistics(env)
116
117        # Instead of starting games immediately for each episodes
118        # The env sample a few random number of "no-op" as a way to introduce randomness.
119        env = NoopResetEnv(env, noop_max=30)
120
121        # Return only every skip-th frame (frameskipping)
122        # And return the max between the two last frames.
123        # (https://stable-baselines3.readthedocs.io/en/master/common/atari\_wrappers.html#stable-baselines3.common.atari\_wrappers.MaxAndSkipEnv)
124        env = MaxAndSkipEnv(env, skip=4)
125
126        # Make end-of-life == end-of-episode, but only reset on true game over.
127        # (https://stable-baselines3.readthedocs.io/en/master/common/atari\_wrappers.html#stable-baselines3.common.atari\_wrappers.EpisodicLifeEnv)
128        env = EpisodicLifeEnv(env)
129
130        # Used for Atari environments that remain static until a "FIRE" action is taken.
131        # Without this, the environment remains static until the agent takes FIRE actions.
132        # As a result, the agent may spend many timestep at the beginning with a static environment until it fires.
133        # Which lead to wasting time steps.
134        # With this, the environment reset as the agent fires. Therefore, the agent are placed in meaningful states without wasting
135        # (https://stable-baselines3.readthedocs.io/en/master/common/atari\_wrappers.html#stable-baselines3.common.atari\_wrappers.FireResetEnv)
136        if "FIRE" in env.unwrapped.get_action_meanings():
137            env = FireResetEnv(env)
138
139        # Clip the reward to {+1, 0, -1} by its sign.
140        # Simplifying rewards helps the agent focus on the direction of improvement (good/bad/neutral) rather than the exact magnitude
141        # (https://stable-baselines3.readthedocs.io/en/master/common/atari\_wrappers.html#stable-baselines3.common.atari\_wrappers.ClipRewardEnv)
142        env = ClipRewardEnv(env)
143
144        # Reduce state information - Resize observations to a smaller size
145        env = gym.wrappers.ResizeObservation(env, (84, 84))
146
147        # Reduce state information - Convert observation into grayscale
148        env = gym.wrappers.GrayScaleObservation(env)
149
150        # Return 4 frame as 1 state - allowing overcoming temporal limitations
151        # Essentially means that the observation or state that it will be returned in the shape of (4, H, W)
152        # 4 because each frame is grayscale, and only has a channel of 1
153        env = gym.wrappers.FrameStack(env, 4)
154
155        # Set the seed for the action space
156        env.action_space.seed(seed)
157
158        return env
159
160    return thunk
161
162
163 # ALGO LOGIC: initialize agent here:
164 class QNetwork(nn.Module):

```

```

165 """
166 A Q-Learning Network.
167
168 Args:
169     env (gym.Env): An environment.
170 """
171 def __init__(self, env):
172     super().__init__()
173     self.network = nn.Sequential(
174         nn.Conv2d(4, 32, 8, stride=4),
175         nn.ReLU(),
176         nn.Conv2d(32, 64, 4, stride=2),
177         nn.ReLU(),
178         nn.Conv2d(64, 64, 3, stride=1),
179         nn.ReLU(),
180         nn.Flatten(),
181         nn.Linear(3136, 512),
182         nn.ReLU(),
183         nn.Linear(512, env.single_action_space.n),
184     )
185
186 def forward(self, x):
187     """
188     Forward propagation for the Q-Learning Network.
189
190     Args:
191         x (float tensor): An observation or state, expected in the shape of (B,C,H,W), in the range of 0 - 255
192
193     Returns:
194         out (float tensor): Probability of taking each actions, in the shape of (B, n) where n is the number of unique actions,
195     """
196     return self.network(x / 255.0)
197
198
199 def linear_schedule(start_e: float, end_e: float, duration: int, t: int):
200     """
201     A linear scheduler that reduces the decaying of epsilon for controlling exploration / exploitation.
202
203     Args:
204         start_e (float): Starting epsilon. \n
205         end_e (float): Ending epsilon. \n
206         duration (int): Total number of episodes.
207         t (int): Current number of episodes.
208
209     Returns:
210         epsilon (float): The epsilon for the current number of episode.
211     """
212     slope = (end_e - start_e) / duration
213     return max(slope * t + start_e, end_e)
214
215
216
217 if __name__ == "__main__":
218     import stable_baselines3 as sb3
219
220     if sb3.__version__ < "2.0":
221         raise ValueError(
222             """Ongoing migration: run the following command to install the new dependencies:
223
224 poetry run pip install "stable_baselines3==2.0.0a1" "gymnasium[atari,accept-rom-license]==0.28.1" "ale-py==0.8.1"
225 """
226         )
227
228     #####
229     # Obtain the arguments with the help of tyro #
230     #####
231     args = tyro.cli(Args)
232
233     # Only allow 1 environment
234     assert args.num_envs == 1, "vectorized envs are not supported at the moment"
235
236     # Create run names
237     run_name = f"{args.env_id}_{args.exp_name}_{args.seed}_{int(time.time())}"
238
239     # Setting up wandb
240     if args.track:
241         import wandb
242
243         wandb.init(
244             project=args.wandb_project_name,
245             entity=args.wandb_entity,
246             sync_tensorboard=True,
247             config=vars(args)

```

```

247     config_vars=args,
248     name=run_name,
249     monitor_gym=True,
250     save_code=True,
251 )
252 writer = SummaryWriter(f"runs/{run_name}")
253 writer.add_text(
254     "hyperparameters",
255     "|param|value|\n|-|\n%s" % ("\n".join([f"|{key}|{value}|" for key, value in vars(args).items()])),
256 )
257
258 # TRY NOT TO MODIFY: seeding
259 random.seed(args.seed)
260 np.random.seed(args.seed)
261 torch.manual_seed(args.seed)
262 torch.backends.cudnn.deterministic = args.torch_deterministic
263
264 device = torch.device("cuda" if torch.cuda.is_available() and args.cuda else "cpu")
265
266 # env setup
267 # With SyncVectorEnv - the environment reset once truncated or terminated
268 # This was done under the hood and not need explicitly coded.
269 envs = gym.vector.SyncVectorEnv(
270     [make_env(args.env_id, args.seed + i, i, args.capture_video, run_name) for i in range(args.num_envs)]
271 )
272
273 # Check if the action_space are discrete
274 assert isinstance(envs.single_action_space, gym.spaces.Discrete), "only discrete action space is supported"
275
276 # Create q_network
277 q_network = QNetwork(envs).to(device)
278
279 # Create optimizer
280 optimizer = optim.Adam(q_network.parameters(), lr=args.learning_rate)
281
282 # Create a copy of the q_network with the same weight and bias
283 target_network = QNetwork(envs).to(device)
284 target_network.load_state_dict(q_network.state_dict())
285
286 # Create a Replay Buffer
287 rb = ReplayBuffer(
288     args.buffer_size,
289     envs.single_observation_space,
290     envs.single_action_space,
291     device,
292     optimize_memory_usage=True,
293     handle_timeout_termination=False,
294 )
295 start_time = time.time()
296
297 # TRY NOT TO MODIFY: start the game
298 obs, _ = envs.reset(seed=args.seed)
299 for global_step in range(args.total_timesteps):
300     # ALGO LOGIC: put action logic here
301
302     # Update epsilon
303     epsilon = linear_schedule(args.start_e, args.end_e, args.exploration_fraction * args.total_timesteps, global_step)
304
305     # Sample an action based on epsilon
306     # Exploration
307     if random.random() < epsilon:
308         actions = np.array([envs.single_action_space.sample() for _ in range(envs.num_envs)])
309
310     # Exploitation
311     else:
312         q_values = q_network(torch.Tensor(obs).to(device))
313         actions = torch.argmax(q_values, dim=1).cpu().numpy()
314
315     # TRY NOT TO MODIFY: execute the game and log data.
316     # Take a step in that direction
317     next_obs, rewards, terminations, truncations, infos = envs.step(actions)
318
319     # TRY NOT TO MODIFY: record rewards for plotting purposes
320     # for wandb
321     if "final_info" in infos:
322         for info in infos["final_info"]:
323             if info and "episode" in info:
324                 print(f"global_step={global_step}, episodic_return={info['episode']['r']}")
325                 writer.add_scalar("charts/episodic_return", info["episode"]["r"], global_step)
326                 writer.add_scalar("charts/episodic_length", info["episode"]["l"], global_step)
327
328     # TRY NOT TO MODIFY: save data to replay buffer; handle `final_observation`
329     real_next_obs = next_obs.copy()

```

```

330
331     # When the episode ends with truncation, the next_obs will not be the real next observation [I assume its because its reset
332     # Instead, it will be stored in infos["final_observation"]
333     # The reason why infos['final_observation'][idx] is used is because it is assuming a vectorised environment
334
335     # No need to perform for the termination state
336     # Because it is still in a meaningful state even after termination (this is before reset)
337
338     # In short
339     # terminated -> in a meaningful state (basically means at the destination) -> then only get reset
340     # truncated -> in a random state (not really meaningful, like step into a wall) -> then only get reset
341     # Therefore, only for truncation that we need to find out and replace the true and meaningful next state
342     for idx, trunc in enumerate(truncations):
343         if trunc:
344             real_next_obs[idx] = infos["final_observation"][idx]
345
346     # Add the information to the replay buffer
347     # When adding these information into the buffer,
348     # For instance if obs has the shape of (num_envs, 4, 32, 32), it will be broken into (4, 32, 32) x num_envs
349     # Then later on, when called from replay buffer, it will return (Batch_size, 4, 32, 32)
350     rb.add(obs, real_next_obs, actions, rewards, terminations, infos)
351
352     # TRY NOT TO MODIFY: CRUCIAL step easy to overlook
353     # ignore the case if terminaed or truncated, i am assuming due to SyncVecEnv, it gets reset in other ways later on
354     obs = next_obs
355
356     # ALGO LOGIC: training.
357     # First check if the number of step is already more than the learning start steps
358     if global_step > args.learning_starts:
359
360         # Check if the step are divisibile by the frequency to decide if gradient descent is performed
361         if global_step % args.train_frequency == 0:
362
363             # Sample data up to the batch size
364             data = rb.sample(args.batch_size)
365
366             # Use the target network to find the next actions based on the recorded next observation / state
367             with torch.no_grad():
368
369                 # returns 2 variables because of using .max
370                 # data.next_observations is in the shape of (B, C, H, W) where B is the number of batch
371                 # target_network(data.next_observations) return (B, n)
372                 # target_network(data.next_observations).max(dim=1) returns the target_max (B,) and its corresponding indices (B, n)
373                 target_max, _ = target_network(data.next_observations).max(dim=1)
374
375                 # Calculatet the td target
376                 # td_target = r + gamma * max(target_network(s', a'), dim = a')
377                 # data.dones is the termination flag
378                 # if termination - data.dones = 1, hence (1 - data.dones.flatten()) = 1 - 1 = 0
379                 # There's no need to predict the action in next state
380                 td_target = data.rewards.flatten() + args.gamma * target_max * (1 - data.dones.flatten())
381
382                 # Calculate the current old value
383                 # old_val = q_network(s, a)
384                 # Which mean it passes the observation of this sample
385                 # q_network(data.observations) -> Send it to the q_network to get (B, n)
386                 # .gather(1, data.actions) -> Based on the index of of the actions, sample the value in the dim=1 (https://docs.pytorch.org/en/stable/tensors.html#torch.Tensor.gather)
387                 # .squeeze -> to reduce the shape to (B,)
388                 old_val = q_network(data.observations).gather(1, data.actions).squeeze()
389
390                 # Compute MSE Loss
391                 loss = F.mse_loss(td_target, old_val)
392
393                 if global_step % 100 == 0:
394                     writer.add_scalar("losses/td_loss", loss, global_step)
395                     writer.add_scalar("losses/q_values", old_val.mean().item(), global_step)
396                     print("SPS:", int(global_step / (time.time() - start_time)))
397                     writer.add_scalar("charts/SPS", int(global_step / (time.time() - start_time)), global_step)
398
399                 # optimize the model
400                 optimizer.zero_grad()
401                 loss.backward()
402                 optimizer.step()
403
404                 # update target network
405                 # Update partially with the help of tau
406                 if global_step % args.target_network_frequency == 0:
407                     for target_network_param, q_network_param in zip(target_network.parameters(), q_network.parameters()):
408                         target_network_param.data.copy_(
409                             args.tau * q_network_param.data + (1.0 - args.tau) * target_network_param.data
410                         )
411
412     if args.save_model:

```

```
413     model_path = f"runs/{run_name}/{args.exp_name}.cleanrl_model"
414     torch.save(q_network.state_dict(), model_path)
415     print(f"model saved to {model_path}")
416     from cleanrl_utils.evals.dqn_eval import evaluate
417
418     episodic_returns = evaluate(
419         model_path,
420         make_env,
421         args.env_id,
422         eval_episodes=10,
423         run_name=f"{run_name}-eval",
424         Model=QNetwork,
425         device=device,
426         epsilon=0.05,
427     )
428     for idx, episodic_return in enumerate(episodic_returns):
429         writer.add_scalar("eval/episodic_return", episodic_return, idx)
430
431     if args.upload_model:
432         from cleanrl_utils.huggingface import push_to_hub
433
434         repo_name = f"{args.env_id}-{args.exp_name}-seed{args.seed}"
435         repo_id = f"{args.hf_entity}/{repo_name}" if args.hf_entity else repo_name
436         push_to_hub(args, episodic_returns, repo_id, "DQN", f"runs/{run_name}", f"videos/{run_name}-eval")
437
438     envs.close()
439     writer.close()
```

 [Show hidden output](#)