

Deep Reinforcement Learning

Homework 2

DQN implementation for Super Mario Bros.

Department: ISA PhD Program

Student ID: 112065802

Name: 蔡睿翊

Preface

In this assignment, I implemented Double DQN to play Super Mario Bros. The environment settings, as well as the processes of training and testing, are introduced in detail as follows.

Training

By following the instructions in the tutorial of Pytorch and modifying the code in the tutorial for performance enhancement, the implementation of training in this assignment includes the following parts.

Environment Preprocessing

Since there are some segments that are not related to the game states in the frames, we use the preprocessing modules including `SkipFrame`, `GrayScaleObservation`, `ResizeObservation`, and `FrameStack` (the last one is imported from `gym.wrappers`) provided by OpenAI Gym to make the representation of observed states concise.

```

22 # Environment preprocessing
23 class SkipFrame(gym.Wrapper):
24     def __init__(self, env, skip):
25         """Return only every `skip`-th frame"""
26         super().__init__(env)
27         self._skip = skip
28
29     def step(self, action):
30         """Repeat action, and sum reward"""
31         total_reward = 0.0
32         for i in range(self._skip):
33             # Accumulate reward and repeat the same action
34             obs, reward, done, info = self.env.step(action)
35             total_reward += reward
36             if done:
37                 break
38         return obs, total_reward, done, info

```

```

40 class GrayScaleObservation(gym.ObservationWrapper):
41     def __init__(self, env):
42         super().__init__(env)
43         obs_shape = self.observation_space.shape[2:]
44         self.observation_space = Box(low=0, high=255, shape=obs_shape, dtype=np.uint8)
45
46     def permute_orientation(self, observation):
47         # permute [H, W, C] array to [C, H, W] tensor
48         observation = np.transpose(observation, (2, 0, 1))
49         observation = torch.tensor(observation.copy(), dtype=torch.float)
50         return observation
51
52     def observation(self, observation):
53         observation = self.permute_orientation(observation)
54         transform = T.Grayscale()
55         observation = transform(observation)
56         return observation

```

```

58 class ResizeObservation(gym.ObservationWrapper):
59     def __init__(self, env, shape):
60         super().__init__(env)
61         if isinstance(shape, int):
62             self.shape = (shape, shape)
63         else:
64             self.shape = tuple(shape)
65
66         obs_shape = self.shape + self.observation_space.shape[2:]
67         self.observation_space = Box(low=0, high=255, shape=obs_shape, dtype=np.uint8)
68
69     def observation(self, observation):
70         transforms = T.Compose([T.Resize(self.shape, antialias=True), T.Normalize(0, 255)])
71         observation = transforms(observation).squeeze(0)
72         return observation

```

Before the training process starts, we apply the aforementioned wrappers to the environment to make the observed states in 4 channels, with each channel having frames resized as 84 by 84 pixels.

```

368         # Apply Wrappers to environment
369         env = SkipFrame(env, skip=4)
370         env = GrayScaleObservation(env)
371         env = ResizeObservation(env, shape=84)
372         env = FrameStack(env, num_stack=4)

```

Double DQN Architecture

The architecture of the dueling DQN is implemented as follows. The architecture includes online and target networks. Each of the networks is built with 3 convolution layers (with ReLU layers in between) followed by the flatten operation and 2 FC layers (with a ReLU layer between them).

```

74 # DQN Architecture
75 class D2QN(nn.Module):
76     """mini CNN structure: input -> (conv2d + relu) x 3 -> flatten -> (dense + relu) x 2 -> output"""
77
78     def __init__(self, input_dim, output_dim):
79         super().__init__()
80         c, h, w = input_dim
81
82         if h != 84:
83             raise ValueError(f"Expecting input height: 84, got: {h}")
84         if w != 84:
85             raise ValueError(f"Expecting input width: 84, got: {w}")
86
87         self.online = self.__build_cnn(c, output_dim)
88
89         self.target = self.__build_cnn(c, output_dim)
90         self.target.load_state_dict(self.online.state_dict())
91
92         # Q_target parameters are frozen.
93         for p in self.target.parameters():
94             p.requires_grad = False
95
96     def forward(self, input, model):
97         if model == "online":
98             return self.online(input)
99         elif model == "target":
100             return self.target(input)
101
102     def __build_cnn(self, c, output_dim):
103         return nn.Sequential(
104             nn.Conv2d(in_channels=c, out_channels=32, kernel_size=8, stride=4),
105             nn.ReLU(),
106             nn.Conv2d(in_channels=32, out_channels=64, kernel_size=4, stride=2),
107             nn.ReLU(),
108             nn.Conv2d(in_channels=64, out_channels=64, kernel_size=3, stride=1),
109             nn.ReLU(),
110             nn.Flatten(),
111             nn.Linear(3136, 512),
112             nn.ReLU(),
113             nn.Linear(512, output_dim),
114         )

```

Agent

The agent is implemented as follows.

Parameters

- Memory size: 100000
- Gamma: 0.9
- Batch size: 32
- Exploration decay in epsilon greedy: 0.95
- Learning rate: 0.0025

```
116 # Agent
117 class Agent:
118     def __init__(self, state_dim, action_dim):
119         self.state_dim = state_dim
120         self.action_dim = action_dim
121
122         self.use_cuda = torch.cuda.is_available()
123
124         # Mario's DNN to predict the most optimal action - we implement this in the Learn section
125         self.net = D2QN(self.state_dim, self.action_dim).float()
126         if self.use_cuda:
127             self.net = self.net.to(device='cuda')
128
129         self.exploration_rate = 1
130         self.exploration_rate_decay = 0.95
131         self.exploration_rate_min = 0.1
132         self.curr_step = 0
133
134         self.save_every = 5e5 # no. of experiences between saving Mario Net
135         self.memory = deque(maxlen=100000)
136         self.batch_size = 32
137
138         self.gamma = 0.9
139
140         self.optimizer = torch.optim.Adam(self.net.parameters(), lr=0.00025)
141         self.loss_fn = torch.nn.SmoothL1Loss()
142
143         self.burnin = 1e4 # min. experiences before training
144         self.learn_every = 3 # no. of experiences between updates to Q online
145         self.sync_every = 1e4 # no. of experiences between Q_target & Q_online sync
146
```

act

This function chooses an action via epsilon-greedy strategy with an exploration decay.

```
147 def act(self, state):
148     """Given a state, choose an epsilon-greedy action"""
149     # EXPLORE
150     if np.random.rand() < self.exploration_rate:
151         action_idx = np.random.randint(self.action_dim)
152
153     # EXPLOIT
154     else:
155         state = state[0].__array__() if isinstance(state, tuple) else state.__array__()
156         state = torch.tensor(state).cuda() if self.use_cuda else torch.tensor(state).unsqueeze(0)
157         action_values = self.net(state, model="online")
158         action_idx = torch.argmax(action_values, axis=1).item()
159
160     # decrease exploration_rate
161     self.exploration_rate *= self.exploration_rate_decay
162     self.exploration_rate = max(self.exploration_rate_min, self.exploration_rate)
163
164     # increment step
165     self.curr_step += 1
166     return action_idx
167
```

cache

This function adds the experiences (states, actions, rewards and corresponding next states) to the memory.

```
168 def cache(self, state, next_state, action, reward, done):
169     """Add the experience to memory"""
170     def first_if_tuple(x):
171         return x[0] if isinstance(x, tuple) else x
172     state = first_if_tuple(state).__array__()
173     next_state = first_if_tuple(next_state).__array__()
174
175     state = torch.FloatTensor(state).cuda() if self.use_cuda else torch.FloatTensor(state)
176     next_state = torch.FloatTensor(next_state).cuda() if self.use_cuda else torch.FloatTensor(next_state)
177     action = torch.LongTensor([action]).cuda() if self.use_cuda else torch.LongTensor([action])
178     reward = torch.DoubleTensor([reward]).cuda() if self.use_cuda else torch.DoubleTensor([reward])
179     done = torch.BoolTensor([done]).cuda() if self.use_cuda else torch.BoolTensor([done])
180
181     self.memory.append((state, next_state, action, reward, done))
```

recall

This function samples a batch of experiences (states, actions, rewards and corresponding next states) from the memory.

```
184 def recall(self):
185     """Sample experiences from memory"""
186     batch = random.sample(self.memory, self.batch_size)
187     state, next_state, action, reward, done = map(torch.stack, zip(*batch))
188     return state, next_state, action.squeeze(), reward.squeeze(), done.squeeze()
```

learn

This function first samples a batch of experiences (states, actions, rewards and corresponding next states) from the memory with function `recall`, estimates the current Q-value with function `td_estimate`, calculates the Q-value in the next state with function `td_target`, and updates the online network with function `update_Q_online` by calculating the TD error between the current Q-value and the Q-value in the next state.

Besides, this function synchronizes the target network and the online network with function `sync_Q_target` per 10000 steps and saves the model with function `save_model` per 500000 steps.

```

190 def td_estimate(self, state, action):
191     current_Q = self.net(state, model="online")[np.arange(0, self.batch_size), action] # Q_online(s,a)
192     return current_Q
193
194 @torch.no_grad()
195 def td_target(self, reward, next_state, done):
196     next_state_Q = self.net(next_state, model="online")
197     best_action = torch.argmax(next_state_Q, axis=1)
198     next_Q = self.net(next_state, model="target")[np.arange(0, self.batch_size), best_action]
199     return (reward + (1 - done.float()) * self.gamma * next_Q).float()
200
201 def update_Q_online(self, td_estimate, td_target):
202     loss = self.loss_fn(td_estimate, td_target)
203     self.optimizer.zero_grad()
204     loss.backward()
205     self.optimizer.step()
206     return loss.item()
207
208 def sync_Q_target(self):
209     self.net.target.load_state_dict(self.net.online.state_dict())
210
211 def learn(self):
212     """Update online action value (Q) function with a batch of experiences"""
213     if self.curr_step % self.sync_every == 0:
214         self.sync_Q_target()
215
216     if self.curr_step % self.save_every == 0:
217         self.save_model()
218
219     if self.curr_step < self.burnin:
220         return None, None
221
222     if self.curr_step % self.learn_every != 0:
223         return None, None
224
225     # Sample from memory
226     state, next_state, action, reward, done = self.recall()
227
228     # Get TD Estimate
229     td_est = self.td_estimate(state, action)
230
231     # Get TD Target
232     td_tgt = self.td_target(reward, next_state, done)
233
234     # Backpropagate loss through Q_online
235     loss = self.update_Q_online(td_est, td_tgt)
236
237     return (td_est.mean().item(), loss)

```

The training process is implemented with 10000 episodes as follows. The metric logger is implemented to monitor the training process (Note the details of the metric logger is omitted here since the main focus of this report is to explain the implementation of double DQN algorithm).

```

380 mario = Agent(state_dim=(4, 84, 84), action_dim=env.action_space.n)
381
382 logger = MetricLogger(save_dir)
383
384 episodes = 10000
385 for e in range(episodes):
386
387     state = env.reset()
388
389     # Play the game!
390     while True:
391
392         # Run agent on the state
393         action = mario.act(state)
394
395         # Agent performs action
396         next_state, reward, done, info = env.step(action)
397
398         # Remember
399         mario.cache(state, next_state, action, reward, done)
400
401         # Learn
402         q, loss = mario.learn()
403
404         # Logging
405         logger.log_step(reward, loss, q)
406
407         # Update state
408         state = next_state
409
410         # Check if end of game
411         if done or info["flag_get"]:
412             break
413
414     logger.log_episode()
415
416     if (e % 20 == 0) or (e == episodes - 1):
417         logger.record(episode=e, epsilon=mario.exploration_rate, step=mario.curr_step)
418         mario.save_model()

```

Testing

Since the trained model is to be loaded, the agent is implemented in a different way from the training process as follows.

- The observed states are preprocessed manually with functions of OpenCV in the function `act`.
- The best action with the maximum probability predicted by the loaded model is chosen in each iteration.

```

65 # Agent
66 class Agent:
67     def __init__(self):
68         # self.use_cuda = torch.cuda.is_available()
69         # Mario's DNN to predict the most optimal action - we implement this in the Learn section
70         self.net = D2QN((4, 84, 84), 12).float()
71         self.load('112065802_hw2_data')
72         # self.net.load_state_dict(torch.load('112065802_hw2_data'), strict=False)
73         self.frames = deque(maxlen=4)
74         self.curr_step = 0
75         self.memory = deque(maxlen=100000)
76
77     def act(self, observation):
78         preprocess_obs = cv2.cvtColor(observation, cv2.COLOR_RGB2GRAY)
79         preprocess_obs = cv2.resize(preprocess_obs, (84, 84), interpolation=cv2.INTER_AREA)
80         while len(self.frames) < 3:
81             self.frames.append(preprocess_obs)
82         self.frames.append(preprocess_obs)
83         preprocess_obs = torch.from_numpy(np.array(self.frames) / 255).float().unsqueeze(0)
84         # observation = observation[0].__array__() if isinstance(observation, tuple) else observation.__array__()
85         # observation = torch.from_numpy(observation.copy()).unsqueeze(0)
86         # print(f'shape of observation: {observation.shape}')
87         action_values = self.net(preprocess_obs, model="online")
88         action_idx = torch.argmax(action_values, axis=1).item()
89
90         # increment step
91         self.curr_step += 1
92
93         return action_idx
94
95     def cache(self, state, next_state, action, reward, done):
96         """Add the experience to memory"""
97         def first_if_tuple(x):
98             return x[0] if isinstance(x, tuple) else x
99         state = first_if_tuple(state).__array__()
100         next_state = first_if_tuple(next_state).__array__()
101
102         state = torch.FloatTensor(state)
103         next_state = torch.FloatTensor(next_state)
104         action = torch.LongTensor([action])
105         reward = torch.DoubleTensor([reward])
106         done = torch.BoolTensor([done])
107
108         self.memory.append((state, next_state, action, reward, done))
109
110     def load(self, load_path):
111         ckp = torch.load(load_path)
112         exploration_rate = ckp.get('exploration_rate')
113         state_dict = ckp.get('model')
114
115         print(f"Loading model at {load_path} with exploration rate {exploration_rate}")
116         self.net.load_state_dict(state_dict)
117         self.exploration_rate = exploration_rate

```

The testing process is implemented with 50 episodes as follows (with each episode iteratively performing the action, cache, and calculation of the cumulative reward). The average reward is calculated after finishing 50 episodes.


```
120 if __name__ == '__main__':
121     env = gym_super_mario_bros.make('SuperMarioBros-v0')
122     env = JoypadSpace(env, COMPLEX_MOVEMENT)
123
124     env.reset()
125     mario = Agent()
126
127     total_reward = 0
128     episodes = 50
129
130     for e in range(episodes):
131         state = env.reset()
132         episode_reward = 0
133         print(f'Episode {e}')
134         while True:
135             # env.render()
136             action = mario.act(state)
137             next_state, reward, done, info = env.step(action)
138             mario.cache(state, next_state, action, reward, done)
139             episode_reward += reward
140
141             state = next_state
142
143             if done or info['flag_get']:
144                 break
145
146         print(f'Episode reward in episode {e}: {episode_reward}')
147         total_reward += episode_reward
148
149     avg_reward = total_reward/50
150     print(f'Average reward: {avg_reward}')
```