# 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 dueling 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.

```python
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
```

```python
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
```

```python
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 Dueling DQN Architecture

The architecture of the dueling double DQN is implemented as follows. The architecture includes online and target networks. Each of the networks is built by a dueling DQN, which consists of 3 convolution layers operated by ReLU, followed by a FC layer for 2 linear layers to predict the value and advantages (for actions in the action space), respectively. Then the two dueling DQNs are applied to construct double dueling DQN as follows.

```python
22      # Dueling Double DQN Architecture
23      class D3QN(nn.Module):
24          def __init__(self, input_dim, output_dim):
25              super().__init__()
26              c, h, w = input_dim
27
28              if h != 84:
29                  raise ValueError(f"Expecting input height: 84, got: {h}")
30              if w != 84:
31                  raise ValueError(f"Expecting input width: 84, got: {w}")
32
33              self.conv1 = nn.Conv2d(c, 4, 3, padding=1)
34              self.conv2 = nn.Conv2d(4, 8, 3, padding=1)
35              self.conv3 = nn.Conv2d(8, 16, 3, padding=1)
36              self.conv4 = nn.Conv2d(16, 16, 3, padding=1)
37              self.conv5 = nn.Conv2d(16, 16, 3, padding=1)
38
39              self.pool = nn.MaxPool2d(2, ceil_mode=True)
40
41              self.fcval = nn.Linear(144, 20)
42              self.fcval2 = nn.Linear(20, 1)
43              self.fcadv = nn.Linear(144, 20)
44              self.fcadv2 = nn.Linear(20, output_dim)
45
46          def forward(self, x):
47              x = self.pool(F.relu(self.conv1(x)))
48              x = self.pool(F.relu(self.conv2(x)))
49              x = self.pool(F.relu(self.conv3(x)))
50              x = self.pool(F.relu(self.conv4(x)))
51              x = self.pool(F.relu(self.conv5(x)))
52
53              x = x.reshape(x.shape[0], -1)
54
55              advantage = F.relu(self.fcadv(x))
56              advantage = self.fcadv2(advantage)
57              advantage = advantage - torch.mean(advantage, dim=-1, keepdim=True)
58
59              value = F.relu(self.fcval(x))
60              value = self.fcval2(value)
61              return value, advantage
```

# 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

```python
114    # Agent
115    class Agent:
116        def __init__(self, state_dim, action_dim):
117            self.state_dim = state_dim
118            self.action_dim = action_dim
119
120            self.use_cuda = torch.cuda.is_available()
121
122            # Mario's DNN to predict the most optimal action - we implement this in the Learn section
123            self.online_net = D3QN(self.state_dim, self.action_dim).float()
124            self.target_net = D3QN(self.state_dim, self.action_dim).float()
125            if self.use_cuda:
126                self.online_net = self.online_net.to(device='cuda')
127                self.target_net = self.target_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.online_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
```

## act

This function chooses an action via an 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().unsqueeze(0) 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 the value and advantage of the current state predicted by the online network in function `td_estimate`, calculates the Q-value in the next state with the values and advantages predicted by online and target networks in function `td_target`, and updates the online network with SGD in 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.

```python
190    def td_estimate(self, state, action):
191        current_V, current_A = self.online_net(state)
192        current_Q = current_V + (current_A - current_A.mean(dim=1, keepdim=True)) # q - batch_size * n_actions
193        return current_Q[np.arange(0, self.batch_size), action]  # Q_online(s,a)
194
195    @torch.no_grad()
196    def td_target(self, reward, next_state, done):
197        next_state_V1, next_state_A1 = self.online_net(next_state)
198        next_state_V2, next_state_A2 = self.target_net(next_state)
199        next_state_Q1 = next_state_V1 + (next_state_A1 - next_state_A1.mean(dim=1, keepdim=True))
200        next_state_Q2 = next_state_V2 + (next_state_A2 - next_state_A2.mean(dim=1, keepdim=True))
201        best_action = torch.argmax(next_state_Q1, axis=1)
202        next_Q2 = next_state_Q2[np.arange(0, self.batch_size), best_action]
203        return (reward + (1 - done.float()) * self.gamma * next_Q2).float()
204
205    def update_Q_online(self, td_estimate, td_target):
206        loss = self.loss_fn(td_estimate, td_target)
207        self.optimizer.zero_grad()
208        loss.backward()
209        self.optimizer.step()
210        return loss.item()
211
212    def sync_Q_target(self):
213        self.target_net.load_state_dict(self.online_net.state_dict())
214
215    def learn(self):
216        """Update online action value (Q) function with a batch of experiences"""
217        if self.curr_step % self.sync_every == 0:
218            self.sync_Q_target()
219
220        if self.curr_step % self.save_every == 0:
221            self.save_model()
222
223        if self.curr_step < self.burnin:
224            return None, None
225
226        if self.curr_step % self.learn_every != 0:
227            return None, None
228
229        # Sample from memory
230        state, next_state, action, reward, done = self.recall()
231
232        # Get TD Estimate
233        td_est = self.td_estimate(state, action)
234
235        # Get TD Target
236        td_tgt = self.td_target(reward, next_state, done)
237
238        # Backpropagate loss through Q_online
239        loss = self.update_Q_online(td_est, td_tgt)
240
241        return (td_est.mean().item(), loss)
```

Then the online network model is saved with the function `save_model`. (PS. Thanks to the TA's suggestion, I found the cause of the model's oversize and updated the way of model saving.)

```python
243    def save_model(self):
244        save_path = '112065802_hw2_data'
245        torch.save(self.online_net.state_dict(), save_path)
246        print(f"D3QN model saved to {save_path} at step {self.curr_step}")
```

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).

```python
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.

```python
61    # Agent
62    class Agent:
63        def __init__(self):
64            # for local test
65            self.use_cuda = torch.cuda.is_available()
66            self.net = D3QN((4, 84, 84), 12)
67            if self.use_cuda:
68                self.net.load_state_dict(torch.load('112065802_hw2_data')).cpu()
69            else:
70                self.net.load_state_dict(torch.load('112065802_hw2_data'))
71            self.frames = deque(maxlen=4)
72            self.curr_step = 0
73            self.memory = deque(maxlen=100000)
74
75        def act(self, observation):
76            preprocess_obs = cv2.cvtColor(observation, cv2.COLOR_RGB2GRAY)
77            preprocess_obs = cv2.resize(preprocess_obs, (84, 84), interpolation=cv2.INTER_AREA)
78            while len(self.frames) < 3:
79                self.frames.append(preprocess_obs)
80            self.frames.append(preprocess_obs)
81            preprocess_obs = torch.from_numpy(np.array(self.frames) / 255).float().unsqueeze(0)
82            _, action_values = self.net(preprocess_obs)
83            action_idx = torch.argmax(action_values, axis=1).item()
84
85            # increment step
86            self.curr_step += 1
87
88            return action_idx
89
90        def cache(self, state, next_state, action, reward, done):
91            """Add the experience to memory"""
92            def first_if_tuple(x):
93                return x[0] if isinstance(x, tuple) else x
94            state = first_if_tuple(state).__array__()
95            next_state = first_if_tuple(next_state).__array__()
96
97            state = torch.FloatTensor(state.copy())
98            next_state = torch.FloatTensor(next_state.copy())
99            action = torch.LongTensor([action])
100           reward = torch.DoubleTensor([reward])
101           done = torch.BoolTensor([done])
102
103           self.memory.append((state, next_state, action, reward, done))
```

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.

```python
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}')
```

# Postscript

Since I cannot use GPU due to the issues of the driver in my lab's servers until I found another resource on 4/10 (Wednesday), and there was a mistake found in the testing code (but fixed at 1am on 4/11) when uploading my trained model to leaderboard, the progress was delayed again……(I also handed in my homework 1 one day late)

However, to take responsibility for my own homework, I decided to update and resubmit the training/testing code (with Double Dueling DQN implementation and model saving/loading), the saved (online) model (with size reduced to 55 KB), and this report.