# 2024 AI - HW#5 Reinforcement Learning

生醫電資所 碩一 **r12945040** 郭思言

## ♦ Describe the Deep Q-Network

The DQN implementation uses a convolutional neural network to estimate Q-values for actions based on image input states. The agent interacts with the environment, stores experiences in a replay buffer, and periodically samples from this buffer to learn. The target network helps stabilize training by providing consistent Q-value estimates.

### ❖ Initialization `__init__`:

```python
 68  class DQN:
 69
 70      def __init__(
 71          self,
 72          state_dim,
 73          action_dim,
 74          lr=1e-4,
 75          epsilon=0.9,
 76          epsilon_min=0.05,
 77          gamma=0.99,
 78          batch_size=64,
 79          warmup_steps=5000,
 80          buffer_size=int(1e5),
 81          target_update_interval=10000,
 82      ):
 83          """
 84          DQN agent has four methods.
 85
 86          - __init__() as usual
 87          - act() takes as input one state of np.ndarray and output actions by following epsilon-greedy policy.
 88          - process() method takes one transition as input and define what the agent do for each step.
 89          - learn() method samples a mini-batch from replay buffer and train q-network
 90          """
 91          self.action_dim = action_dim
 92          self.epsilon = epsilon
 93          self.gamma = gamma
 94          self.batch_size = batch_size
 95          self.warmup_steps = warmup_steps
 96          self.target_update_interval = target_update_interval
 97          self.network = PacmanActionCNN(state_dim[0], action_dim)
 98          self.target_network = PacmanActionCNN(state_dim[0], action_dim)
 99          self.target_network.load_state_dict(self.network.state_dict())
100          self.optimizer = torch.optim.Adam(self.network.parameters(), lr)
101
102          self.buffer = ReplayBuffer(state_dim, (1,), buffer_size)
103          self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
104          self.network.to(self.device)
105          self.target_network.to(self.device)
106
107          self.total_steps = 0
108          self.epsilon_decay = (epsilon - epsilon_min) / 1e6
```

1. Initializes the CNN model for action selection (`network`) and its target network (`target_network`).
2. Sets hyperparameters such as learning rate, epsilon for exploration, discount factor (gamma), batch size, etc.
3. Initializes the replay buffer and optimizer.
4. Sets device for computation (CPU/GPU) and initializes other variables like epsilon decay and total steps.

❖ Method `act`:

```
110    @torch.no_grad()
111    def act(self, x, training=True):
112        self.network.train(training)
113        if training and ((np.random.rand() < self.epsilon) or (self.total_steps < self.warmup_steps)):
114            # Random action
115            action = np.random.randint(0, self.action_dim)
116        else:
117            # output actions by following epsilon-greedy policy
118            x = torch.from_numpy(x).float().unsqueeze(0).to(self.device)
119
120            # get q-values from network
121            q_value = self.network(x)
122
123            # get action with maximum q-value
124            action = np.argmax(q_value.cpu().data.numpy())
125        return action
```

1. Selects an action based on the epsilon-greedy policy.

2. Chooses a random action with probability epsilon, otherwise selects the action with the highest Q-value.

❖ Method `learn`:

```
127    def learn(self):
128        # sample a mini-batch from replay buffer
129        state, action, reward, next_state, terminated = map(
130            lambda x: x.to(self.device), self.buffer.sample(self.batch_size)
131        )
132
133        # get q-values from network
134        q_values = self.network(state)  # shape: [64, 9]
135        next_q_values = self.target_network(next_state).detach()  # shape: [64, 9]
136
137        # compute td_target
138        # td_target: if terminated, only reward, otherwise reward + gamma * max(next_q)
139        td_target = reward + (1 - terminated) * self.gamma * next_q_values.max(1)[0].view(-1, 1)  # shape: [64, 1]
140
141        # gather the Q-values for the actions that were actually taken
142        q_value = q_values.gather(1, action.long())  # use long to convert to int
143
144        # compute loss with td_target and q-values
145        loss = F.mse_loss(q_value, td_target)
146
147        # initialize optimizer
148        self.optimizer.zero_grad()
149        # backpropagation
150        loss.backward()
151        nn.utils.clip_grad_norm_(self.network.parameters(), 10)
152
153        # update network
154        self.optimizer.step()
155        return {"value_loss": loss.item()}  # return dictionary for logging
```

1. Samples a batch of experiences from the replay buffer.

2. Computes Q-values for the current states and next states using the network and target network, respectively.

3. Calculates the TD target, considering whether the episode has terminated.

4. Computes the loss between the Q-values of the taken actions and the TD target.

5. Performs backpropagation to update the network's weights.

❖ Method `process`:

```
157        def process(self, transition):
158
159            result = {"value_loss": 0}
160            self.total_steps += 1
161
162            # update replay buffer
163            # transition: (state, action, reward, next_state, terminated)
164            self.buffer.update(*transition)
165
166            if self.total_steps > self.warmup_steps:
167                result = self.learn()
168
169            if self.total_steps % self.target_update_interval == 0:
170                # update target network
171                "self.target_network.YOUR_CODE_HERE"
172                self.target_network.load_state_dict(self.network.state_dict())
173
174            self.epsilon -= self.epsilon_decay
175            return result
176
```
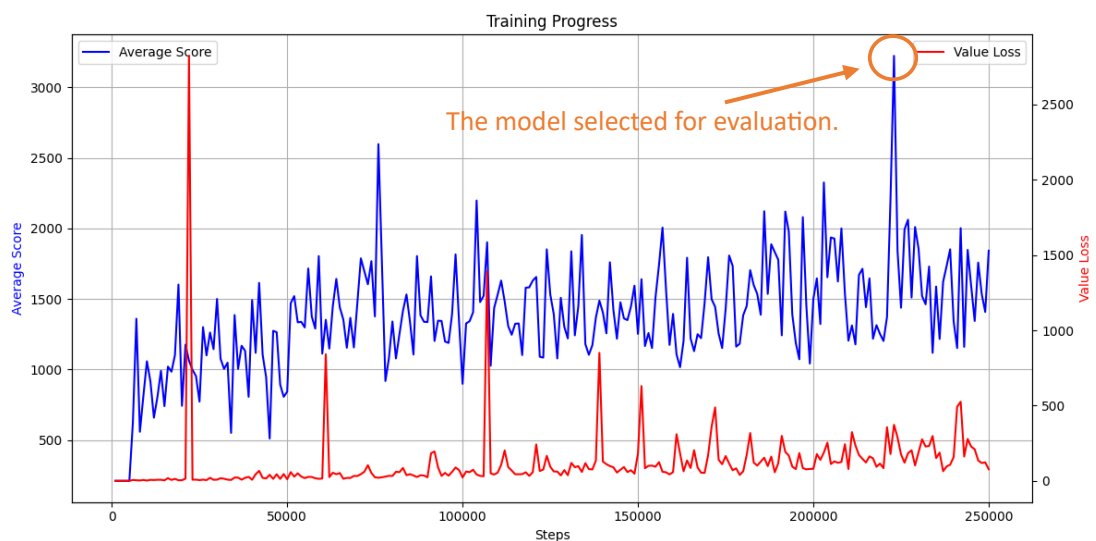
1. Processes a single transition, updating the replay buffer and performing learning if sufficient steps have been taken.
2. Updates the target network at regular intervals and decays epsilon for exploration.


♦ **Describe the architecture of your PacmanActionCNN**

```
10    class PacmanActionCNN(nn.Module):
11
12        def __init__(self, state_dim, action_dim):
13            super(PacmanActionCNN, self).__init__()
14            # build your own CNN model
15            self.conv1 = nn.Conv2d(state_dim, 32, kernel_size=6, stride=2)
16            self.bn1 = nn.BatchNorm2d(32)
17            self.conv2 = nn.Conv2d(32, 64, kernel_size=4, stride=2)
18            self.bn2 = nn.BatchNorm2d(64)
19
20            self.hidden1 = nn.Linear(64 * 19 * 19, 1024)
21            self.hidden2 = nn.Linear(1024, 128)
22            self.output = nn.Linear(128, action_dim)
23
24        def forward(self, x):
25            x = F.relu(self.bn1(self.conv1(x)))
26            x = F.relu(self.bn2(self.conv2(x)))
27            # print(x.shape)
28            x = torch.flatten(x, start_dim=1)
29            x = F.relu(self.hidden1(x))
30            x = F.relu(self.hidden2(x))
31            return self.output(x)
```

❖ My PacmanActionCNN contains 2 2D Convolutional layers, both accompanied by Batch Normalization and ReLU activation. The first 2D Convolutional layer has a kernel size of 6 and a stride of 2, while the second 2D Convolutional layer has a kernel size of 4 and a stride of 2. These layers extract features from the input state, scaling the dimensions to 64x19x19.

❖ Following the convolutional layers, I apply a Multi-Layer Perceptron (MLP) with 2 hidden layers to transform the features into the final output dimension.

♦ **Plot your training curve, including both loss and rewards**



♦ **Show screenshots from your evaluation video**

**("ALE/MsPacman-v5 has a total of three chances. Display the reward (score) each time you are caught.)**