

Part 1-1: Minimax Search (10%)

```
140 # Begin your code
141 def minimax_search(gameState, depth, agentIndex):
142     # terminal states
143     if gameState.isWin() or gameState.isLose() or depth == self.depth:
144         return self.evaluationFunction(gameState)
145
146     # Pacman's turn
147     if agentIndex == 0:
148         maxScore = -float("inf")
149         actions = gameState.getLegalActions(0)
150
151         # get maximum
152         for action in actions:
153             nextState = gameState.getNextState(0, action)
154             maxScore = max(
155                 maxScore, minimax_search(nextState, depth, 1))
156
157         return maxScore
158
159     # ghosts' turns
160     else:
161         minScore = float("inf")
162         actions = gameState.getLegalActions(agentIndex)
163
164         # get minimum
165         for action in actions:
166             nextState = gameState.getNextState(agentIndex, action)
167
168             # next agent is the Pacman
169             if agentIndex == (gameState.getNumAgents() - 1):
170                 minScore = min(minScore, minimax_search(
171                     nextState, depth + 1, 0))
172
173             # next agent is a ghost
174             else:
175                 minScore = min(minScore, minimax_search(
176                     nextState, depth, agentIndex + 1))
177
178         return minScore
179
180     # initial condition
181     actions = gameState.getLegalActions(0)
182     maxScore = -float("inf")
183     returnAction = None
184
185     # get maximum and decide next action
186     for action in actions:
187         nextState = gameState.getNextState(0, action)
188         score = minimax_search(nextState, 0, 1)
189
190         if score > maxScore:
191             returnAction = action
192             maxScore = score
193
194     return returnAction
195 # End your code
```

In the function `minimax_search()`, there are three parameters: `gameState`, `depth` and `agentIndex`. While `gameState` represents the current condition of the game, we record the depth of the search tree to check if we have gone through enough predictions. Every time when `minimax_search()` is called, we first check if it's in the terminal condition (win or lose), or we have reached the bottom of the search tree (by checking depth). Then, if `agentIndex = 0`, which means the Pacman should act, we perform a depth first search on all legal actions to find the maximum score of the particular action. On the other hand, if `agentIndex > 0`, which means a ghost should act, we also perform a depth first search on all legal actions to find the minimum score of the particular action. Notice that in ghosts' turns, we should check if the next agent is the Pacman or a ghost to ensure we transfer to the right state. Finally, set the initial condition (the Pacman starts first), and choose the action which will lead to maximum score gained.

```

*** Running MinimaxAgent on smallClassic 1 time(s).
Pacman died! Score: 84
Average Score: 84.0
Scores:      84.0
Win Rate:    0/1 (0.00)
Record:      Loss
*** Finished running MinimaxAgent on smallClassic after 0 seconds.
*** Won 0 out of 1 games. Average score: 84.000000 ***
*** PASS: test_cases\part1-1\8-pacman-game.test

### Question part1-1: 10/10 ###

```

Part 1-2: Expectimax Search (10%)

```

198 class ExpectimaxAgent(MultiAgentSearchAgent):
199     """
200     | Your expectimax agent (part1-2)
201     """
202
203     def getAction(self, gameState):
204         """
205         Returns the expectimax action using self.depth and self.evaluationFunction
206
207         All ghosts should be modeled as choosing uniformly at random from their
208         legal moves.
209         """
210         """ YOUR CODE HERE """
211         # Begin your code
212         def expectimax_search(gameState, depth, agentIndex):
213             # terminal states
214             if gameState.isWin() or gameState.isLose() or depth == self.depth:
215                 return self.evaluationFunction(gameState)
216
217             # Pacman's turn
218             if agentIndex == 0:
219                 maxScore = -float("inf")
220                 actions = gameState.getLegalActions(0)
221
222                 # get maximum
223                 for action in actions:
224                     nextState = gameState.getNextState(0, action)
225                     maxScore = max(
226                         maxScore, expectimax_search(nextState, depth, 1))

```

```

227
228         return maxScore
229
230     # ghosts' turns
231     else:
232         actions = gameState.getLegalActions(agentIndex)
233         averageScore = 0
234
235         # get expected value (mean)
236         for action in actions:
237             nextState = gameState.getNextState(agentIndex, action)
238
239             # next agent is the Pacman
240             if agentIndex == (gameState.getNumAgents() - 1):
241                 averageScore += (expectimax_search(nextState,
242                                                     depth + 1, 0) / len(actions))
243
244             # next agent is a ghost
245             else:
246                 averageScore += (expectimax_search(nextState,
247                                                     depth, agentIndex + 1) / len(actions))
248
249         return averageScore
250

```

```

251     # initial condition
252     actions = gameState.getLegalActions(0)
253     maxScore = -float("inf")
254     returnAction = None
255
256     # get maximum and decide next action
257     for action in actions:
258         nextState = gameState.getNextState(0, action)
259         score = expectimax_search(nextState, 0, 1)
260
261         if score > maxScore:
262             returnAction = action
263             maxScore = score
264
265     return returnAction
266     # End your code
267
268
269     better = scoreEvaluationFunction
270

```

The only difference between part 1-1 and part 1-2 is that in `expectimax_search()`, we don't expect the ghosts to choose the worst actions towards us (Pacman). In contrast, we assume the actions of ghosts are random (with equal probability). Therefore, we just need to change the formula (line 240 ~ 247) of the returned score in ghosts' turns (`agentIndex > 0`).

```

*** Running ExpectimaxAgent on smallClassic 1 time(s).
Pacman died! Score: 84
Average Score: 84.0
Scores:      84.0
Win Rate:    0/1 (0.00)
Record:      Loss
*** Finished running ExpectimaxAgent on smallClassic after 0 seconds.
*** Won 0 out of 1 games. Average score: 84.000000 ***
*** PASS: test_cases\part1-2\7-pacman-game.test

### Question part1-2: 10/10 ###

```

Part 2-1: Value Iteration (10%)

```

70     def runValueIteration(self):
71         # Write value iteration code here
72         """ YOUR CODE HERE """
73         # Begin your code
74         for i in range(self.iterations):
75             tempValues = util.Counter()
76             states = self.mdp.getStates()
77
78             for state in states:
79                 if not self.mdp.isTerminal(state):
80                     Q_values = util.Counter()
81                     actions = self.mdp.getPossibleActions(state)
82
83                     for action in actions:
84                         Q_values[action] = self.computeQValueFromValues(
85                             state, action)
86
87                     tempValues[state] = max(Q_values.values())
88
89             self.values = tempValues
90         # End your code

```

In each iteration (the outer loop), we go through all states and all actions to find the maximum Q value for each state (by calling `computeQValueFromValues()`). Keep the value in a table (`tempValues`), and update `self.values` in every iteration.

```
98 def computeQValueFromValues(self, state, action):
99     """
100     Compute the Q-value of action in state from the
101     value function stored in self.values.
102     """
103     """ YOUR CODE HERE """
104     # Begin your code
105     QValue = 0
106     transitionStatesAndProbabilities = self.mdp.getTransitionStatesAndProbs(
107         state, action)
108
109     for (nextState, probability) in transitionStatesAndProbabilities:
110         QValue += probability * \
111             (self.mdp.getReward(state, action, nextState) +
112              self.discount * self.values[nextState])
113
114     return QValue
115     # End your code
```

In `computeQValueFromValues()`, we calculate Q value, which is the expected total reward, by summing the total reward of each action times the probability to make that action. The total reward involves current reward (`self.mdp.getReward(state, action, nextState)`) and the future reward, which will be multiplied by a discount factor (`self.discount * self.values[nextState]`).

```
117 def computeActionFromValues(self, state):
118     """
119     The policy is the best action in the given state
120     according to the values currently stored in self.values.
121
122     You may break ties any way you see fit. Note that if
123     there are no legal actions, which is the case at the
124     terminal state, you should return None.
125     """
126     """ YOUR CODE HERE """
127     # Begin your code
128
129     # check for terminal
130     if self.mdp.isTerminal(state):
131         return None
132
133     else:
134         QValues = util.Counter()
135         actions = self.mdp.getPossibleActions(state)
136
137         for action in actions:
138             QValues[action] = self.computeQValueFromValues(state, action)
139
140         return QValues.argmax()
141     # End your code
```

If the current state is a terminal state, then return no action (`None`). Otherwise, return the action which is corresponding with the largest Q value.

Question part2-1

=====

```
*** PASS: test_cases\part2-1\1-tinygrid.test
*** PASS: test_cases\part2-1\2-tinygrid-noisy.test
*** PASS: test_cases\part2-1\3-bridge.test
*** PASS: test_cases\part2-1\4-discountgrid.test

### Question part2-1: 10/10 ###
```

Part 2-2: Q-learning (10%)

```
49     def __init__(self, **args):
50         "You can initialize Q-values here..."
51         ReinforcementAgent.__init__(self, **args)
52
53         "**** YOUR CODE HERE ****"
54         # Begin your code
55         self.values = util.Counter()
56         # End your code
57
58     def getQValue(self, state, action):
59         """
60         Returns Q(state,action)
61         Should return 0.0 if we have never seen a state
62         or the Q node value otherwise
63         """
64         "**** YOUR CODE HERE ****"
65         # Begin your code
66         return self.values[(state, action)]
67         # End your code
```

It would return 0.0 if we have never seen a state since `self.values = util.Counter()` is initialized with 0.0.

```
69     def computeValueFromQValues(self, state):
70         """
71         Returns max_action Q(state,action)
72         where the max is over legal actions. Note that if
73         there are no legal actions, which is the case at the
74         terminal state, you should return a value of 0.0.
75         """
76         "**** YOUR CODE HERE ****"
77         # Begin your code
78         actions = self.getLegalActions(state)
79
80         return 0.0 if not actions else max([self.getQValue(state, action) for action in actions])
81         # End your code
```

If no legal actions, return 0. Otherwise, return the maximum Q value.

```

83     def computeActionFromQValues(self, state):
84         """
85         Compute the best action to take in a state. Note that if there
86         are no legal actions, which is the case at the terminal state,
87         you should return None.
88         """
89         """ YOUR CODE HERE """
90         # Begin your code
91         actions = self.getLegalActions(state)
92         maxQValue = self.computeValueFromQValues(state)
93
94         return None if not actions else random.choice([action for action in actions if self.getQValue(state, action) == maxQValue])
95         # End your code

```

As comments.

```

119     def update(self, state, action, nextState, reward):
120         """
121         The parent class calls this to observe a
122         state = action => nextState and reward transition.
123         You should do your Q-Value update here
124
125         NOTE: You should never call this function,
126         it will be called on your behalf
127         """
128         """ YOUR CODE HERE """
129         # Begin your code
130         self.values[(state, action)] = (1 - self.alpha) * \
131             self.values[(state, action)] + self.alpha * (reward +
132                                                         self.discount * self.computeValueFromQValues(nextState))
133         # End your code

```

By the formula:

$$q_{\pi}(s, a) = (1 - \alpha)q_{\pi}(s, a) + \alpha[R + \gamma \max_{a'} q_{\pi}(s', a')]$$

```

Question part2-2
=====

*** PASS: test_cases\part2-2\1-tinygrid.test
*** PASS: test_cases\part2-2\2-tinygrid-noisy.test
*** PASS: test_cases\part2-2\3-bridge.test
*** PASS: test_cases\part2-2\4-discountgrid.test

### Question part2-2: 10/10 ###

```

Part 2-3: epsilon-greedy action selection (5%)

```

97     def getAction(self, state):
98         """
99         Compute the action to take in the current state. With
100         probability self.epsilon, we should take a random action and
101         take the best policy action otherwise. Note that if there are
102         no legal actions, which is the case at the terminal state, you
103         should choose None as the action.
104
105         HINT: You might want to use util.flipCoin(prob)
106         HINT: To pick randomly from a list, use random.choice(list)
107         """
108         # Pick Action
109         legalActions = self.getLegalActions(state)
110         action = None
111         """ YOUR CODE HERE """
112         # Begin your code
113         if util.flipCoin(self.epsilon):
114             return random.choice(legalActions)
115         else:
116             return self.computeActionFromQValues(state)
117         # End your code

```

• You can also observe the following simulations for different epsilon values. Does that behavior of the agent match what you expect?

Answer: `python gridworld.py -a q -k 100 --noise 0.0 -e 0.9`



The behavior is as I expected ($\uparrow \uparrow \rightarrow \rightarrow \rightarrow$). The average returns from the start state will decrease as epsilon becomes large, due to the random actions when $\epsilon > 0$.

Question part2-3

=====

```
*** PASS: test_cases\part2-3\1-tinygrid.test
*** PASS: test_cases\part2-3\2-tinygrid-noisy.test
*** PASS: test_cases\part2-3\3-bridge.test
*** PASS: test_cases\part2-3\4-discountgrid.test

### Question part2-3: 5/5 ###
```

Part 2-4: Approximate Q-learning (10%)

```
196 def getQValue(self, state, action):
197     """
198     Should return Q(state,action) = w * featureVector
199     where * is the dotProduct operator
200     """
201     """ YOUR CODE HERE """
202     # Begin your code
203     # get weights and feature
204     return self.getWeights() * self.featsExtractor.getFeatures(state, action)
205     # End your code
```

As comments.

```
207 def update(self, state, action, nextState, reward):
208     """
209     Should update your weights based on transition
210     """
211     """ YOUR CODE HERE """
212     # Begin your code
213     features = self.featsExtractor.getFeatures(state, action)
214     difference = reward + self.discount * \
215         self.computeValueFromQValues(
216             nextState) - self.getQValue(state, action)
217
218     for feature in features:
219         self.weights[feature] += self.alpha * \
220             difference * features[feature]
221     # End your code
```

By the formula:

$$\begin{aligned} w_i &\leftarrow w_i + \alpha[\text{correction}] f_i(s, a) \\ \text{correction} &= (R(s, a) + \gamma V(s')) - Q(s, a) \end{aligned}$$

Part 3 : DQN (10%)

```
Average Score: 1546.6
Scores:         1756.0, 1159.0, 1558.0, 1711.0, 1549.0
Win Rate:       5/5 (1.00)
Record:         Win, Win, Win, Win, Win
```

```
1  import numpy as np
2  import torch
3  import torch.nn as nn
4  import torch.nn.functional as F
5
6  """ Deep Q Network """
7
8
9  class DQN(nn.Module):
10     def __init__(self, num_inputs=6, num_actions=4):
11         super(DQN, self).__init__()
12
13         self.conv1 = nn.Conv2d(num_inputs, 32, kernel_size=3, stride=1)
14         self.conv2 = nn.Conv2d(32, 64, kernel_size=2, stride=1)
15         self.fc3 = nn.Linear(4352, 512)
16         self.fc4 = nn.Linear(512, num_actions)
17
18     def forward(self, x):
19         x = F.relu(self.conv1(x))
20         x = F.relu(self.conv2(x))
21         # print(x.view(x.size(0), -1).shape)
22         x = F.relu(self.fc3(x.view(x.size(0), -1)))
23         return self.fc4(x)
24
```

```
25 # model parameters
26 model_trained = True
27
28 GAMMA = 0.8 # discount factor
29 LR = 0.02 # learning rate
30
31 batch_size = 32 # memory replay batch size
32 memory_size = 100000 # memory replay size
33 start_training = 5000 # start training at this episode
34 TARGET_REPLACE_ITER = 100 # update network step
35
36 epsilon_final = 0.1 # epsilon final
37 epsilon_step = 10000
```

- What is the difference between On-policy and Off-policy?

Answer: An off-policy learner learns the value of the optimal policy independently of the agent's actions. Q-learning is an off-policy learner. An on-policy learner learns the value of the policy being carried out by the agent including the exploration steps.

- Briefly explain value-based, policy-based and Actor-Critic. Also, describe the value function $V^\pi(S)$.

Answer: In value-based methods, the agent learns to estimate the value of different actions or states in the environment. In policy-based methods, the agent learns directly the policy that maps states to actions. Actor-Critic methods combine the benefits of both value-based and policy-based methods. The agent has two components: an actor

that learns the policy, and a critic that estimates the value function. $V^\pi(S)$ is the expected cumulative reward that an agent can achieve starting from state S and following policy π .

- What is the difference between Monte-Carlo (MC) based approach and Temporal-difference (TD) approach for estimating $V^\pi(S)$?

Answer: The main difference between MC and TD methods is the way they estimate the value function. MC methods rely on the empirical average of returns, which are the cumulative rewards obtained from a particular state or action until the end of the episode. On the other hand, TD methods estimate the value function by bootstrapping, which means updating the estimate using a new estimate. TD methods update the value function incrementally after each time step, by combining the reward obtained in that step with the estimated value of the next state. To be specific, MC methods require the agent to complete an entire episode before updating the value function, while TD methods can update the value function at each time step.

- Describe State-action value function $Q^\pi(s,a)$ and the relationship between $V^\pi(S)$ in Q-learning.

Answer: The state-action value function $Q^\pi(s,a)$ is the expected cumulative reward that an agent can achieve starting from state s , taking action a , and following policy π . The relationship between $Q^\pi(s,a)$ and $V^\pi(S)$ in Q-learning:

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) * Q^\pi(s, a)$$

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V^\pi(s')]$$

- Describe following tips Target Network, Exploration and Replay Buffer using in Q-learning.

Answer: The target network is a copy of the Q-network, and its parameters are frozen for a certain number of iterations. This technique reduces the correlation between the target value and the Q-value and leads to a more stable and efficient learning process. Exploration allows the agent to discover new and potentially better policies. One common exploration strategy used in Q-learning is the epsilon-greedy policy. The epsilon parameter determines the trade-off between exploration and exploitation, and it is gradually reduced over time as the agent learns. Replay buffer is a memory buffer that stores the agent's experience in the form of (state, action, reward, next state) tuples. During learning, the agent samples a batch of experiences from the replay buffer and uses them to update the Q-function.

- Explain what is different between DQN and Q-learning.

Answer: The difference between DQN and Q-learning include their function approximation, experience replay, and the target network. In Q-learning, the Q-values are updated based on the maximum Q-value of the next state after every step based on the current experience. On the other hand, in DQN, it introduces the concept of target network and experience replay, which help to make the learning process more stable, and it uses a neural network to approximate the Q-function, which allows DQN to handle high-dimensional state spaces and continuous action spaces.

- Compare the performance of every method and do some discussions in your report.

Answer:

```
PS D:\courses\Intro-to-AI\HW3\Adversarial_search> python pacman.py -p MinimaxAgent -l trappedClassic -a depth=3 -q -n 10
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Pacman died! Score: -501
Average Score: -501.0
Scores: -501.0, -501.0, -501.0, -501.0, -501.0, -501.0, -501.0, -501.0, -501.0, -501.0
Win Rate: 0/10 (0.00)
Record: Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss
PS D:\courses\Intro-to-AI\HW3\Adversarial_search>
```

```

PS D:\courses\Intro-to-AI\HW3\Adversarial_search> python pacman.py -p ExpectimaxAgent -l trappedClassic -a depth=3 -q -n 10
Pacman emerges victorious! Score: 532
Pacman emerges victorious! Score: 532
Pacman emerges victorious! Score: 532
Pacman died! Score: -502
Pacman emerges victorious! Score: 532
Pacman died! Score: -502
Pacman died! Score: -502
Pacman died! Score: -502
Pacman emerges victorious! Score: 532
Pacman died! Score: -502
Average Score: 15.0
Scores: 532.0, 532.0, 532.0, -502.0, 532.0, -502.0, -502.0, -502.0, 532.0, -502.0
Win Rate: 5/10 (0.50)
Record: Win, Win, Win, Loss, Win, Loss, Loss, Loss, Win, Loss
PS D:\courses\Intro-to-AI\HW3\Adversarial_search>

```

For part 1, by minimax search, Pacman rushes to the closest ghost since in its prediction (worst case), it'll all eventually lose the game. However, it's not the case that the ghosts will choose the worst actions towards Pacman. Therefore, by expectimax search, Pacman will try to grab food and sometimes it wins.

```

-----
Pacman died! Score: -404
Pacman died! Score: -383
Pacman died! Score: -374
Pacman died! Score: -228
Pacman died! Score: -420
Pacman died! Score: -375
Pacman died! Score: -400
Pacman died! Score: -413
Pacman died! Score: -428
Pacman died! Score: -412
Average Score: -383.7
Scores: -404.0, -383.0, -374.0, -228.0, -420.0, -375.0, -400.0, -413.0, -428.0, -412.0
Win Rate: 0/10 (0.00)
Record: Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss

```

(python pacman.py -p ApproximateQAgent -x 2000 -n 2010 -l smallGrid)

```

-----
Pacman emerges victorious! Score: 527
Pacman emerges victorious! Score: 525
Pacman emerges victorious! Score: 529
Pacman emerges victorious! Score: 529
Pacman emerges victorious! Score: 529
Pacman emerges victorious! Score: 527
Pacman emerges victorious! Score: 527
Pacman emerges victorious! Score: 525
Pacman emerges victorious! Score: 527
Pacman emerges victorious! Score: 527
Average Score: 527.2
Scores: 527.0, 525.0, 529.0, 529.0, 529.0, 527.0, 527.0, 525.0, 527.0, 527.0
Win Rate: 10/10 (1.00)
Record: Win, Win, Win, Win, Win, Win, Win, Win, Win, Win

```

(python pacman.py -p ApproximateQAgent -a extractor=SimpleExtractor -x 50 -n 60
-l mediumGrid)

```

PS D:\courses\Intro-to-AI\HW3\Q-learning> python pacman.py -p ApproximateQAgent -a extractor=SimpleExtractor -x 20 -n 25 -l smallClassic
Beginning 20 episodes of Training
Training Done (turning off epsilon and alpha)
-----
Pacman emerges victorious! Score: 955
Pacman died! Score: -127
Pacman died! Score: -291
Pacman emerges victorious! Score: 970
Pacman emerges victorious! Score: 972
Average Score: 495.8
Scores: 955.0, -127.0, -291.0, 970.0, 972.0
Win Rate: 3/5 (0.60)
Record: Win, Loss, Loss, Win, Win

```

(python pacman.py -p ApproximateQAgent -a extractor=SimpleExtractor -x 20 -n 25
-l smallClassic)

```
Average Score: 1484.8
Scores:        1372.0, 1698.0, 1347.0, 1349.0, 1658.0
Win Rate:      5/5 (1.00)
Record:        Win, Win, Win, Win, Win
```

(python pacman.py -p PacmanDQN -n 25 -x 20 -l smallClassic)

For part 2 and part 3, we can see that the default PacmanQAgent performs poorly, while the performance ApproximateQAgent with SimpleExtractor is good. If we put ApproximateQAgent with SimpleExtractor and PacmanDQN in the same situation (same layout and same training times), we can see that the performance of PacmanDQN is even better.

- Describe problems you meet and how you solve them.

Answer: 執行 Part3 時遭遇警告訊息: RuntimeError: Attempting to deserialize object on a CUDA device but torch.cuda.is_available() is False. If you are running on a CPU-only machine, please use torch.load with map_location=torch.device('cpu') to map your storages to the CPU. 解決方法如[連結](#)。此外，寫作業時在實作與公式上有不懂的地方，我主要參考了這些網站：

[MDP 模型之 Grid World\(值迭代方法\) UncoDong 的博客-CSDN 博客](#)

[MDP 模型之 Grid World\(Q Learning 方法\) gridworld 价值迭代 UncoDong 的博客-CSDN 博客](#)

[什么是 Q Learning \(Reinforcement Learning 强化学习\) - YouTube](#)

[李宏毅 DRL Lecture 3: Q-learning \(Basic Idea\) - HackMD](#)