

Part 1

1. Minimax Search

```

1 class MinimaxAgent(MultiAgentSearchAgent):
2     def getAction(self, gameState):
3         """ YOUR CODE HERE """
4         # Begin your code
5         def minimax(agent, depth, gameState):
6             if gameState.isWin() or gameState.isLose():
7                 return self.evaluationFunction(gameState), None
8             if depth == self.depth:
9                 return self.evaluationFunction(gameState), None
10
11             next_agent = agent + 1
12             if next_agent >= gameState.getNumAgents():
13                 next_agent = 0
14                 depth += 1
15             moves = gameState.getLegalActions(agent)
16             if not moves:
17                 return self.evaluationFunction(gameState), None
18
19             results = []
20             for move in moves:
21                 next_state = gameState.getNextState(agent, move)
22                 score, _ = minimax(next_agent, depth, next_state)
23                 results.append((score, move))
24             if agent == 0:
25                 # Pacman's turn, maximize the value
26                 best_value = max(results, key=lambda x: x[0])
27             else:
28                 # Ghost's turn, minimize the value
29                 best_value = min(results, key=lambda x: x[0])
30             return best_value
31
32         _, action = minimax(0, 0, gameState)
33         return action # Return the action component of the tuple
34         # End your code

```

The **getAction** method uses a recursive **minimax** function to evaluate game states. The function checks for terminal states (win or loss) and the maximum search depth as base cases to return the evaluated score. It then determines the next agent's turn and increases depth appropriately after each cycle of agents. For each possible move, it recursively evaluates future game states, collecting scores for each move. Pacman aims to maximize his score, while ghosts aim to minimize it. The algorithm ultimately returns the optimal move for Pacman based on the evaluated outcomes of all possible actions.

2. Expectimax Search

```

1 class ExpectimaxAgent(MultiAgentSearchAgent):
2     def getAction(self, gameState):
3         """ YOUR CODE HERE """
4         # Begin your code
5         def expectimax(agent, depth, gameState):
6             if gameState.isWin() or gameState.isLose() or depth == self.depth
              * gameState.getNumAgents():
7                 return self.evaluationFunction(gameState)
8
9             actions = gameState.getLegalActions(agent)
10            if agent == 0:
11                # Pacman's turn, maximize
12                return max(expectimax((agent + 1) % gameState.getNumAgents(),
              depth + 1, gameState.getNextState(agent, action))
              for action in actions)
13            else:
14                # Ghosts' turn, expected value
15                return sum(expectimax((agent + 1) % gameState.getNumAgents(),
              depth + 1, gameState.getNextState(agent, action))
              for action in actions) / len(actions)
16
17            # Start from Pacman (agent 0) and depth 0
18            bestAction = max(gameState.getLegalActions(0), key=lambda x:
              expectimax(1, 1, gameState.getNextState(0, x)))
19            return bestAction
20
21            # End your code
22

```

The method employs a recursive **expectimax** function to evaluate game states by considering all possible future scenarios within a specified depth. The function checks for terminal states (win or loss) and stops further recursion when it reaches the maximum depth. For each agent's turn, Pacman aims to maximize his expected utility by choosing the move with the highest score, while the ghosts calculate the expected utility across all possible moves, reflecting their average impact on the game state. The best move for Pacman is selected by evaluating the expected outcomes of all legal actions from the initial state.

Part 2

1. Value Iteration

```

1 class ValueIterationAgent(ValueEstimationAgent):
2     def runValueIteration(self):
3         """ YOUR CODE HERE """
4         # Begin your code
5         for _ in range(self.iterations):
6             newValues = util.Counter()
7             for state in self.mdp.getStates():
8                 if self.mdp.isTerminal(state):
9                     continue
10                max_value = float('-inf')
11                for action in self.mdp.getPossibleActions(state):
12                    q_value = self.computeQValueFromValues(state, action)
13                    if q_value > max_value:
14                        max_value = q_value
15                newValues[state] = max_value if max_value != float('-inf')
16                else 0
17            self.values = newValues
18        # End your code
19
20    def computeQValueFromValues(self, state, action):
21        """ YOUR CODE HERE """
22        # Begin your code
23        q_value = 0
24        trans = self.mdp.getTransitionStatesAndProbs(state, action)
25        for nextState, prob in trans:
26            reward = self.mdp.getReward(state, action, nextState)
27            q_value += prob * (reward + self.discount * self.values[nextState])
28        return q_value
29        # End your code
30
31    def computeActionFromValues(self, state):
32        """ YOUR CODE HERE """
33        # Begin your code
34        #check for terminal
35        if self.mdp.isTerminal(state):
36            return None
37        best_action = None
38        max_value = float('-inf')
39        for action in self.mdp.getPossibleActions(state):
40            q_value = self.computeQValueFromValues(state, action)
41            if q_value > max_value:
42                max_value = q_value
43                best_action = action
44        return best_action
45        # End your code

```

The **runValueIteration** method performs the main iteration loop for a specified number of iterations. In each iteration, it updates the value for each non-terminal state based on the maximum expected utility, calculated through the **computeQValueFromValues** method, which computes the Q-value for a given state-action pair by summing the expected utilities of possible future states. The algorithm assumes future rewards are discounted by a factor. Once the iterations are complete, the **computeActionFromValues** method determines the optimal action for any state by selecting the action with the highest Q-value.

2. Q-learning

```

1 class QLearningAgent(ReinforcementAgent):
2     def __init__(self, **args):
3         "You can initialize Q-values here..."
4         ReinforcementAgent.__init__(self, **args)
5         "*** YOUR CODE HERE ***"
6         # Begin your code
7         self.qValues = util.Counter()
8         # End your code
9
10
11     def getQValue(self, state, action):
12         "*** YOUR CODE HERE ***"
13         # Begin your code
14         return self.qValues[(state, action)]
15         # End your code
16
17     def computeValueFromQValues(self, state):
18         "*** YOUR CODE HERE ***"
19         # Begin your code
20         legalActions = self.getLegalActions(state)
21         if not legalActions:
22             return 0.0
23         return max(self.getQValue(state, action) for action in legalActions)
24         # End your code
25
26     def computeActionFromQValues(self, state):
27         "*** YOUR CODE HERE ***"
28         # Begin your code
29         legalActions = self.getLegalActions(state)
30         if not legalActions:
31             return None
32         bestVal = self.computeValueFromQValues(state)
33         bestActions = [action for action in legalActions if self.getQValue(
34             state, action) == bestVal]
35         return random.choice(bestActions)
36         # End your code
37
38     def update(self, state, action, nextState, reward):
39         "*** YOUR CODE HERE ***"
40         # Begin your code
41         sample = reward + self.discount * self.computeValueFromQValues(

```

```

        nextState)
41     self.qValues[(state, action)] = (1 - self.alpha) * self.getQValue(
        state, action) + self.alpha * sample
42     # End your code

```

- Q-value Retrieval

The **getQValue** method retrieves the current Q-value for a given state-action pair, defaulting to zero if it has not been previously encountered.

- Value Calculation

The **computeValueFromQValues** method determines the highest Q-value for any action available from a given state, returning zero if there are no legal actions.

- Action Determination

The **computeActionFromQValues** method identifies the action that yields the highest Q-value for a given state. If multiple actions have the same value, one is randomly selected, providing an exploration mechanism.

- Q-value Update

The **update** method adjusts the Q-values based on new data from the environment. It computes a sample Q-value using the reward received for taking an action and the discounted maximum Q-value of the resulting state. The Q-value for the state-action pair is then updated towards this sample, scaled by a learning rate.

3. Epsilon-greedy Action Selection

```

1  class QLearningAgent(ReinforcementAgent):
2      def getAction(self, state):
3          # Pick Action
4          legalActions = self.getLegalActions(state)
5          action = None
6          """ YOUR CODE HERE """
7          # Begin your code
8          if legalActions:
9              if util.flipCoin(self.epsilon):
10                 action = random.choice(legalActions)
11             else:
12                 action = self.computeActionFromQValues(state)
13         return action
14         # End your code

```

The **getAction** method in the **QLearningAgent** class determines the next action for the agent to take from a given state. It first retrieves all legal actions available from the current state. The method then employs an ϵ -greedy strategy to balance exploration and exploitation:

- Exploration

With a probability of ϵ , the agent randomly selects an action from the available legal actions. This randomness introduces exploration, allowing the agent to discover new strategies.

- Exploitation

With a probability of $1 - \epsilon$, the agent selects the action that maximizes the expected utility based on current Q-values, calculated through the **computeActionFromQValues** method. This approach ensures that the agent exploits its current knowledge to maximize the expected reward.

4. Approximate Q-learning (Bonus)

```

1 class ApproximateQAgent(PacmanQAgent):
2     def getQValue(self, state, action):
3         """ YOUR CODE HERE """
4         # Begin your code
5         # get weights and feature
6         features = self.feateXtractor.getFeatures(state, action)
7         return sum(self.weights[f] * value for f, value in features.items())
8         # End your code
9
10    def update(self, state, action, nextState, reward):
11        """ YOUR CODE HERE """
12        # Begin your code
13        features = self.feateXtractor.getFeatures(state, action)
14        correction = (reward + self.discount * self.computeValueFromQValues(
15            nextState)) - self.getQValue(state, action)
16        for feature, value in features.items():
17            self.weights[feature] += self.alpha * correction * value
18        # End your code

```

- Q-value Computation

The **getQValue** method calculates the estimated Q-value for a given state-action pair. It does this by extracting relevant features of the state-action context through `feateXtractor.getFeatures`, and then computing the dot product of these features with their corresponding weights. This results in a linear approximation of the Q-value based on the weighted sum of the features.

- Q-value Update

For each feature in the current state-action pair, the corresponding weight is updated to better align the predicted Q-values with the observed values. The update rule applies a learning rate **alpha**, the correction term, and the feature's value to adjust the weight, promoting convergence to optimal policies over time.

Results:

```

1 $ python pacman.py -p ApproximateQAgent -a extractor=SimpleExtractor -x 50 -n
   60 -l mediumGrid
2 Average Score: 527.8
3 Scores:      529.0, 527.0, 525.0, 527.0, 529.0, 529.0, 529.0, 525.0, 529.0,
   529.0
4 Win Rate:    10/10 (1.00)
5 Record:      Win, Win, Win, Win, Win, Win, Win, Win, Win, Win

```

Part 3

Results:

```

1 $ python pacman.py -p PacmanDQN -n 25 -x 20 -l smallClassic -q
2 Average Score: 1391.6
3 Scores:          964.0, 1512.0, 1179.0, 1764.0, 1539.0
4 Win Rate:        5/5 (1.00)
5 Record:          Win, Win, Win, Win, Win

```

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

- **On-policy** methods learn the value of the policy being carried out by the agent, including the exploration steps. The policy that is evaluated and improved is the same policy used to make decisions. A classic example of an on-policy method is **SARSA**.

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$$

$$S \leftarrow S', A \leftarrow A'$$

- **Off-policy** methods learn the value of an optimal policy independently of the agent's actions, including learning from actions that are outside the current policy. A prominent example of an off-policy method is **Q-learning**.

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$$S \leftarrow S'$$

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

- Value-based Methods

These methods focus on finding the value function that estimates how good it is for an agent to be in a given state. The policy is derived from the value function, typically by choosing the action that maximizes the expected value. Q-learning and Value Iteration are common examples of value-based methods.

- Policy-based Methods

These methods directly parameterize and optimize the policy that dictates the action to be taken in each state. They do not explicitly compute value functions. Policy Gradient is a typical example where the policy is directly optimized through gradients.

- Actor-Critic Methods

These combine features of both value-based and policy-based methods. The "actor" updates the policy distribution in the direction suggested by the "critic," which evaluates the action taken by the actor according to a value function. This approach aims to use the strengths of both value estimation and policy optimization.

- Value Function $V^\pi(s)$

The value function under a policy π , denoted as $V^\pi(s)$, is the expected return (sum of discounted rewards) starting from state s and following policy π thereafter. Mathematically, it's given by:

$$V^\pi(s) = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, \pi\right]$$

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

- Monte-Carlo Methods

These methods estimate $V^\pi(s)$ by averaging the returns observed after visiting a state s . This approach does not update estimates based on other estimates, and each update requires a complete episode to finish, making it only applicable to episodic tasks.

- Temporal-Difference Methods: TD methods estimate $V^\pi(s)$ by updating estimates based on other learned estimates, without waiting for the episode to end. They update estimates based partly on the observed reward and partly on estimates of future states. This allows TD methods to learn online and from incomplete episodes.

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

The state-action value function $Q^\pi(s, a)$ represents the expected return starting from state s , taking an action a , and thereafter following policy π . The relationship between $Q^\pi(s, a)$ and $V^\pi(s)$ is given by:

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

In Q-learning, which is a value-based method, the optimal state-action value function $Q^*(s, a)$ is learned, and the optimal policy can be derived by selecting the action that maximizes $Q^*(s, a)$ at each state.

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

- Target Network

In Deep Q-learning, a separate network (target network) is used to estimate the target values during updates, which helps in stabilizing the learning. The target network's weights are periodically updated to the weights of the primary network.

- Exploration

This refers to the technique used to balance the exploration of the environment (trying new things) and exploitation (using known information). Common strategies include ϵ -greedy, where ϵ represents the probability of choosing a random action.

- Replay Buffer

This is a technique where experiences (state, action, reward, next state) are stored in a buffer and sampled randomly to update the network. This helps in breaking the correlation between consecutive learning updates and stabilizing learning.

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

- **DQN** extends Q-learning by using a deep neural network to approximate the Q-value function. Unlike traditional Q-learning, which uses a table to store Q-values, DQN can handle environments with high-dimensional state spaces.
- **Q-learning** is a tabular method generally suitable for smaller discrete state spaces and can suffer from scalability issues in complex environments. DQN, on the other hand, can generalize across similar states due to its function approximation approach.

The following table shows a clear progression in performance from simpler to more advanced Pacman agents. The ReflexAgent performs poorly with the lowest score and no wins, while Minimax and Expectimax show moderate improvements, with Expectimax notably achieving a positive average score and a 23% win rate. The advanced ApproximateQAgent and PacmanDQN display the best performance, achieving the highest win rates and scores respectively.

	ReflexAgent	MinimaxAgent	ExpectimaxAgent	ApproximateQAgent	PacmanDQN
AvgScore	-284.25	-127.0	134.46	816.62	1222.83
WinRate	0/100 (0.00)	6/100 (0.06)	23/100 (0.23)	87/100 (0.87)	80/100 (0.80)

[illegible]

```
1 $ python pacman.py -p MinimaxAgent -n 100 -l smallClassic -q
2 Average Score: -127.0
3 Scores:      -196.0, -167.0, -174.0, 67.0, 140.0, -175.0, -329.0, -174.0, -94.0, -448.0, 14.0,
      -216.0, -613.0, 1257.0, -519.0, -315.0, -389.0, -47.0, 1095.0, -372.0, 1077.0, -16.0, 64.0,
      -336.0, -444.0, -33.0, -273.0, -374.0, -192.0, -360.0, -215.0, -260.0, 0.0, 408.0, -219.0,
      -49.0, -240.0, 467.0, 157.0, -202.0, -182.0, -266.0, -175.0, -290.0, -241.0, -73.0, -162.0,
      -495.0, 33.0, -184.0, -161.0, -353.0, -262.0, -140.0, -516.0, 1043.0, 38.0, -184.0, -252.0,
      -206.0, -28.0, -260.0, -220.0, -184.0, -68.0, 47.0, -335.0, -376.0, 198.0, -214.0, -180.0,
      -134.0, -464.0, -105.0, 408.0, -212.0, -310.0, -211.0, -368.0, 148.0, -29.0, -221.0, -386.0,
      -174.0, -399.0, -455.0, -308.0, -128.0, -328.0, -329.0, -159.0, -409.0, -226.0, 770.0,
      23.0, -249.0, -423.0, 16.0, -594.0, -101.0
4 Win Rate:      6/100 (0.06)
5 Record:      Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Win,
      Loss, Loss, Loss, Loss, Win, Loss, Win, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss,
      Loss, Loss, Loss, Loss, Win, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss,
      Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Win, Loss, Loss, Loss,
      Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss,
      Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss, Loss,
      Loss, Loss, Loss, Loss, Win, Loss, Loss, Loss, Loss, Loss, Loss
```


Appendix

1. Reference

- (a) OpenAI. (2024). ChatGPT (4) [Large language model]. <https://chat.openai.com>
- (b) Hung-yi lee. (2021, May 22). 【機器學習 2021】概述增強式學習 (Reinforcement Learning, RL). <https://shorturl.at/inFR4>