



**UNIVERSITATEA  
TEHNICĂ  
DIN CLUJ-NAPOCA**

---

**Probleme de cautare si agenti adversariali**

*Inteligența Artificială*

---

Autori: Istrate Cristian Andrei si Ivan Alexandru Sofron  
Grupa: 30237

FACULTATEA DE AUTOMATICA  
SI CALCULATOARE

16 Ianuarie 2024

## Cuprins

<b>1</b>	<b>Reinforcement Learning . . . . .</b>	<b>2</b>
1.1	Question 1 - Value iteration . . . . .	2
1.2	Question 3 - Policies . . . . .	3
1.3	Question 5 - Q-Learning . . . . .	4

# 1 Reinforcement Learning

## 1.1 Question 1 - Value iteration

Ca prima cerinta am avut de realizat un "Value iteration agent", care este un planificator offline asa ca ceea ce ne intereseaza pe noi este ca in functie de fiecare iteratie pe care o face spre finalul labirintului sa asociem mutarilor cate o valoare in functie de cat de aproape suntem de final. La fiecare final de joc, valorile se actualizeaza in functie de cat de buna a fost mutarea spre finalul jocului.

Cheia spre implementarea acestui algoritm a fost ecuatia Bellman (figura 1). Astfel am avut de implementat 3 functii care au dus la construirea si aplicarea algoritmului Bellman asupra problemei noastre.

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

Figura 1: Ecuatia Bellman

```
1  def runValueIteration(self):
2      # Write value iteration code here
3      "* YOUR CODE HERE *
4      for _ in range(self.iterations):
5          new_values = util.Counter()
6          for state in self.mdp.getStates():
7              if not self.mdp.isTerminal(state):
8                  new_values[state] = max(self.computeQValueFromValues(state, action)
9                      for action in self.mdp.getPossibleActions(state))
10             self.values = new_values
11
12  def computeQValueFromValues(self, state, action):
13      """
14          Compute the Q-value of action in state from the
15          value function stored in self.values.
16          """
17      "* YOUR CODE HERE *
18      q_value = 0
19      for next_state, prob in self.mdp.getTransitionStatesAndProbs(state, action):
20          reward = self.mdp.getReward(state, action, next_state)
21          q_value += prob * (reward + self.discount * self.values[next_state])
22      return q_value
23      # util.raiseNotDefined()
24
25  def computeActionFromValues(self, state):
26      """
27          The policy is the best action in the given state
28          according to the values currently stored in self.values.
29
```

```

30         You may break ties any way you see fit. Note that if
31         there are no legal actions, which is the case at the
32         terminal state, you should return None.
33         """
34         "* YOUR CODE HERE *
35         possible_actions = self.mdp.getPossibleActions(state)
36
37         if not possible_actions:
38             return None
39
40         best_q_value = float('-inf')
41         best_action = None
42
43         for action in possible_actions:
44             q_value = self.computeQValueFromValues(state, action)
45             if q_value > best_q_value:
46                 best_action = action
47                 best_q_value = q_value
48
49         return best_action
50         # util.raiseNotDefined()

```

Pentru a reusi implementarea functiei Bellman a trebuit sa intelegem ecuatia Bellman si sa o impartim in bucati pentru a putea implementa algoritmul de care aveam nevoie.

## 1.2 Question 3 - Policies

Pentru a doua cerinta am avut de stabilit valori pentru discount, noise si living reward astfel incat sa parcurgem labirintul in cele 5 cazuri. Cele 3 valori semnifica:

- answerDiscount: Reprezinta cat de mult ne concentram pe viitoarele recompense, mai mic insemnand ca mergem pe cele mai apropiate recompense, iar mai mare pe cele departate.
- answerNoise: Reprezinta probabilitatea de efectua o miscare aleatoare, diferita de miscarea aleasa de noi.
- answerLivingReward: Reprezinta cat de mult ne dorim sa ramanem intr-o stare, astfel incrementand aceasta variabila, petrecem mai mult timp fara sa incheiem jocul.

```
python gridworld.py -g DiscountGrid -a value --discount [YOUR_DISCOUNT] --noise [YOUR_NOISE] --livingReward [YOUR_LIVING_REWARD]opy
```

Figura 2: GUI command

Dupa multiple teste cu comanda din figura 2 am ajuns la urmatoarele valori:

```

1 def question3a():
2     answerDiscount = 0.1
3     answerNoise = -0.1
4     answerLivingReward = 0.0
5     return answerDiscount, answerNoise, answerLivingReward
6     # If not possible, return 'NOT POSSIBLE'

```

```

7
8 def question3b():
9     answerDiscount = 0.1
10    answerNoise = 0.1
11    answerLivingReward = 0.0
12    return answerDiscount, answerNoise, answerLivingReward
13    # If not possible, return 'NOT POSSIBLE'
14
15 def question3c():
16    answerDiscount = 0.5
17    answerNoise = 0.0
18    answerLivingReward = 0.0
19    return answerDiscount, answerNoise, answerLivingReward
20    # If not possible, return 'NOT POSSIBLE'
21
22 def question3d():
23    answerDiscount = 0.5
24    answerNoise = 0.1
25    answerLivingReward = 0.0
26    return answerDiscount, answerNoise, answerLivingReward
27    # If not possible, return 'NOT POSSIBLE'
28
29 def question3e():
30    answerDiscount = 0.0
31    answerNoise = 0.0
32    answerLivingReward = 1.0
33    return answerDiscount, answerNoise, answerLivingReward
34    # If not possible, return 'NOT POSSIBLE'

```

### 1.3 Question 5 - Q-Learning

Pentru ultima cerinta am avut de completat mai multe metode ale clasei "QLearningAgent". Acest agent nu mai foloseste modelul MDP pe care il folosea "Value iteration" ci se bazeaza pe modelul de invatare prin incercare, astfel se foloseste de mediul in care se afla pentru a-si calcula miscarile. Pentru a reusi implementarea acestor metode ne-am folosit de clasa "ValueIteration" pe care am implementat-o la Q1 pentru a avea un punct de plecare.

Metodele care difera de Q1 sunt: computeValueFromQValues si update. Prima metoda a fost implementata cu ajutorul comentariilor oferite ca suport, iar cea de a doua reprezenta formula de baza pentru a implementa Q learning agent si a fost necesar sa cautam mai multe informatii despre acest agent, intr-un final gasind formula care a dus la rezolvarea functiei (Fig. 3).

$$Q^{new}(S_t, A_t) \leftarrow (1 - \underbrace{\alpha}_{\text{learning rate}}) \cdot \underbrace{Q(S_t, A_t)}_{\text{current value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left( \underbrace{R_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(S_{t+1}, a)}_{\text{estimate of optimal future value}} \right)}_{\text{new value (temporal difference target)}}$$

Figura 3: Q Learning equation

```

1  def _init_(self, **args):
2      "You can initialize Q-values here..."
3      ReinforcementAgent._init_(self, **args)
4      "* YOUR CODE HERE *"
5      self.qValue = util.Counter()
6
7  def getQValue(self, state, action):
8      """
9          Returns Q(state,action)
10         Should return 0.0 if we have never seen a state
11         or the Q node value otherwise
12     """
13     "* YOUR CODE HERE *"
14     return self.qValue[(state, action)]
15     # util.raiseNotDefined()
16
17
18  def computeValueFromQValues(self, state):
19      """
20          Returns max_action Q(state,action)
21          where the max is over legal actions. Note that if
22          there are no legal actions, which is the case at the
23          terminal state, you should return a value of 0.0.
24     """
25     "* YOUR CODE HERE *"
26     legalActions = self.getLegalActions(state)
27     if not legalActions:
28         return 0.0
29     return max(self.getQValue(state, action) for action in legalActions)
30     # util.raiseNotDefined()
31
32  def computeActionFromQValues(self, state):
33      """
34          Compute the best action to take in a state. Note that if there
35          are no legal actions, which is the case at the terminal state,
36          you should return None.
37     """
38     "* YOUR CODE HERE *"
39     legalActions = self.getLegalActions(state)
40     if not legalActions:
41         return None
42
43     maxQValue = float('-inf')
44     bestAction = None
45
46     for action in legalActions:
47         qValue = self.getQValue(state, action)
48         if qValue > maxQValue:

```

```

49         maxQValue = qValue
50         bestAction = action
51
52     return bestAction
53     # util.raiseNotDefined()
54
55 def update(self, state, action, nextState, reward):
56     """
57     The parent class calls this to observe a
58     state = action => nextState and reward transition.
59     You should do your Q-Value update here
60
61     NOTE: You should never call this function,
62     it will be called on your behalf
63     """
64     "* YOUR CODE HERE *"
65     currentQValue = self.getQValue(state, action)
66     nextQValue = self.computeValueFromQValues(nextState)
67
68     updatedQValue = (1 - self.alpha) * currentQValue + self.alpha * (reward +
69     self.discount * nextQValue)
70     self.qValue[(state, action)] = updatedQValue
71     # util.raiseNotDefined()

```