# Reinforcement Learning

*Inteligenta Artificiala*

Autori: Laura-Luisa Voicu
Grupa: 30236

Autori: Laura-Luisa Voicu
Grupa: 30236

FACULTATEA DE AUTOMATICA
SI CALCULATOARE

11 Ianuarie 2024

# Cuprins

# 1 Introducere

## 1.1 Overview

### 1.1.1 Despre Reinforcement Learning

Reinforcement Learning = tehnica de invatare automata prin care un software e invatat sa ia decizii pentru a obtine rezultatele cele mai optime.

Actiunile software-ului ce contribuie la atingerea obiectivului vor fi intarite, in timp ce actiunile ce indeparteaza obiectivul sunt ignorate.

Reinforcement Learning - types:

1. Passive Reinforcement Learning: how to learn from already given experiences
   (a) Model-based: modelul MDP invata din experiente, dupa care rezolva MDP-ul
   (b) Model-free:
      i. Value learning: invata valori dintr-un policy fix
      ii. Q Learning: invata Q valori din policy-ul optimal
2. Active Reinforcement Learning: hot to collect new experiences
3. Approximate Reinforcement Learning: to handle large state spaces

## 1.2 MDP - Markov Decision Process

Reinforcement Learning foloseste MDP
- Are un set de stari
- Are un set de actiuni per stare
- Are un model T(s,a,s')
- Are o functie de recompensa R(s,a,s')

## 1.3 MDPs - Instructiuni

Rulare Gridworld (foloseste arrow keys)

```
python3 gridworld.py -m

# agentii se muta random
python3 gridworld.py -g MazeGrid
```

In Gridworld MDP pozitiile sunt reprezentate ca si coordonate (x,y). Default, recompensa este 0; ea poate fi schimbata cu optiunea -r.

# 2 Question 1: Value Iteration

## 2.1 Update State

Pentru a updata starea curenta, se va folosi formula:

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

Figura 1: Update next state

Q1 presupune scrierea unui value iteration agent ce reprezinta un **offline planner**, nu un agent de reinforcement learning.

Clasa ValueIterationAgent ca un MDP si ruleaza iteratiile de atatea ori cat a fost specificat.

Se vor computa k pasi ce estimeaza valoarea optimala Vk (in metoda **runValueIteration**. E necesara implementarea :

- computeActionFromValue(state) - are rolul de a computa cea mai buna actiune corespunzatoare valorii din self.values
- computeQValueFromValues(state, action) - returneaza Q-value al tuplei (state, action) corespunzatoare valorii din self.values.

## 2.2  Comenzi de rulare

```
python3 gridworld.py -a value -i 100 -k 10
python3 gridworld.py -a value -i 5

#Autograder:
python3 autograder.py -q q1
```
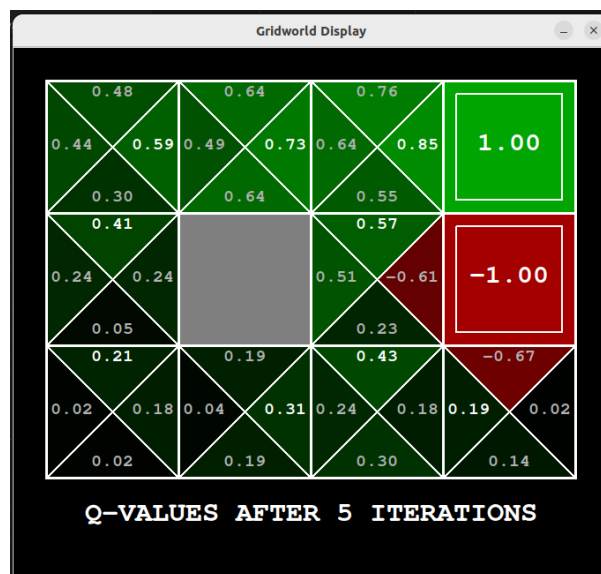


Figura 2: Value Iteration

## 2.3  Cod

:

```
1  class ValueIterationAgent(ValueEstimationAgent):
2
3      def runValueIteration(self):
4          for _ in range(self.iterations):
5              aux_values = util.Counter()
6              for state in self.mdp.getStates():
7                  if not self.mdp.isTerminal(state):
```

3

```
8              actions = self.mdp.getPossibleActions(state)
9              max_value_state = self.computeQValueFromValues(state, actions[0])
10             for action in actions:
11                 value = self.computeQValueFromValues(state, action)
12                 if value > max_value_state:
13                     max_value_state = value
14             aux_values[state] = max_value_state
15           else:
16             aux_values[state] = 0
17       self.values = aux_values
18
19
20   def computeQValueFromValues(self, state, action):
21
22       #V(k+1)(s) = max(Sum(T(s,a,s')*[R(s,a,s') + gama*V(k)(s')]))
23       # s' -  noua stare, V - values, R - reward, gama - factor discount
24       q_value = 0
25       R = self.mdp.getReward
26       gama = self.discount
27       V = self.values
28       all_states_probs = self.mdp.getTransitionStatesAndProbs(state, action)
29       for new_state, prob in all_states_probs:
30         q_value += prob * ( R(state, action, new_state) + gama * V[new_state] )
31       return q_value
32
33
34   def computeActionFromValues(self, state):
35       if self.mdp.isTerminal(state):
36         return None
37
38       # state neterminal
39       policies = [(self.computeQValueFromValues(state, action), action) for action in self
40
41        # returnam cheia pentru care policy e maxim
42       return max(policies, key = lambda x : x[0])[1]
```

### 2.4   Dificultati la implementare

- Intelegerea conceptelor: MDP, Reinforcement Learning, Values / QValues
- Integrarea functilor deja existente prezentate pentru MDP (getTransitionStatesAndProbs, etc.)
- Integrarea efectiva a functiei prezentate mai sus pentru calculul urmatoarei stari.

## 3   Question 3: Policies

### 3.1   Context

**DiscountGrid** - Ce stare terminala alegem?

Exista 2 tipuri de path-uri:

- Risk the cliff (1) - mai scurte dar risca sa castige un payoff negativ mai mare.
- Avoid the cliff (2) - mai lungi dar e mai putin probabil sa gastige payoff negativ mare.
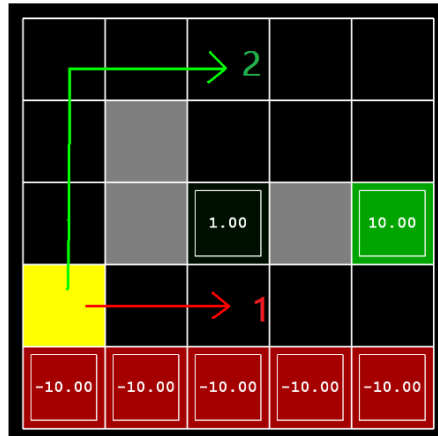


Figura 3: Cliffs

## 3.2    Cerinta

Setarea discount-ului, noise-ului si a reward-ului astfel incat sa se produca policy-ul optimal pentru diverse tipuri.

Optimal policy types

- Prefer the close exit (+1), risking the cliff (-10)
- Prefer the close exit (+1), but avoiding the cliff (-10)
- Prefer the distant exit (+10), risking the cliff (-10)
- Prefer the distant exit (+10), avoiding the cliff (-10)
- Avoid both exits and the cliff (so an episode should never terminate)

## 3.3    Comenzi de rulare

```
python3 gridworld.py -g DiscountGrid -a value --discount [YOUR_DISCOUNT] --noise [YOUR_NOISE
```

```
#Autograder:
python3 autograder.py -q q3
```

## 3.4    Cod

```
1   ef question3a():
2       answerDiscount = 0.5
3       answerNoise = 0
4       answerLivingReward = -1
5       return answerDiscount, answerNoise, answerLivingReward
6
7   def question3b():
8       answerDiscount = 0.5
```

```
9        answerNoise = 0.25
10       answerLivingReward = -1
11       return answerDiscount, answerNoise, answerLivingReward
12
13   def question3c():
14       answerDiscount = 1
15       answerNoise = 0
16       answerLivingReward = -1
17       return answerDiscount, answerNoise, answerLivingReward
18
19   def question3d():
20       answerDiscount = 0.5
21       answerNoise = 0.5
22       answerLivingReward = 0.5
23       return answerDiscount, answerNoise, answerLivingReward
24
25   def question3e():
26       answerDiscount = 1
27       answerNoise = 0
28       answerLivingReward = 1
29       return answerDiscount, answerNoise, answerLivingReward
30
31   def question8():
32       answerEpsilon = None
33       answerLearningRate = None
34       return answerEpsilon, answerLearningRate
35
```

### 3.5   Justificari

- Question3a
  - answerDiscount: 0.5 - accent moderat pe reward imediat
  - answerNoise: 0 - actiuni cunoscute, nu random
  - abswerLivingReward: -1 penalizarea agentului daca continua
  - concluzie: Agentul e focusat pe reward imediat si primeste penalizare daca contiuna. Setarile implica tinderea spre castiguri de scurta durata, dar adoptand o abordare prudenta (answerNoice e 0).
- Question3b
  - answerDiscount: 0.5 - accent moderat pe reward imediat
  - answerNoise: 0.25 - pot aparea decizii random in actiunea agentului
  - abswerLivingReward: -1 - penalizarea agentului daca continua
  - concluzie: Agentul manifesta un echilibru intre recompensele imediate si cele indepartate.
- Question3c
  - answerDiscount: 1 - trateaza in mod egal recompensele imediate si cele indepartate
  - answerNoise: 0 - nu exista actiuni random in actiunile agentului.
  - abswerLivingReward: -1 - penalizarea agentului daca continua

- coneluzie: Agentul acorda aceeasi prioritate recompenselor imediate si indepartate, dar absenta Noise-ului poate duce la abordare prudenta intrucat actiunile viitoare sunt cunoscute.
- Question3d
  - answerDiscount: 1 - trateaza in mod egal recompensele imediate si cele indepartate
  - answerNoise: 0.5 - apar decizii random in actiunea agentului
  - abswerLivingReward: 0.5 - se ofera o recompensa povitiva pentru existenta continua a a gentului
  - coneluzie: Agentul explora ma mult si e mai dispus la riscuri.
- Question3e
  - answerDiscount: 1 - trateaza in mod egal recompensele imediate si cele indepartate
  - answerNoise: 0 - actiuni cunoscute, nu random
  - abswerLivingReward: 1 - se ofera o recompensa povitiva pentru existenta continua a a gentului
  - coneluzie: Agenntul are o existenta prelungita, iar recompensele indepartate si imeditate sunt tratate in mod egal.
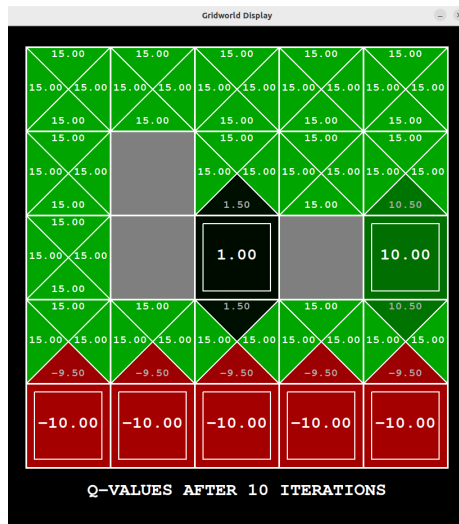


Figura 4: Values

Figura 5: QValues

# 4 Question 5: Q-Learning

## 4.1 Context

Agentul nu invata efectiv din experienta, ci doar din modelul MDP. Cand agentul interactioneaza cu environment-ul, el doar urmeaza policy-urile setate.

Agentul Q-learning invata din incercarile si erorile interactiunilor cu environmen-ul prin metoda **update**.

## 4.2 Comenzi de rulare

```
python3 gridworld.py -a q -k 5 -m
```

```
#Autograder:
python3 autograder.py -q q5
```

## 4.3 Cod

```
1  class QLearningAgent(ReinforcementAgent):
2
3      def __init__(self, **args):
4          self.values = util.Counter()
5
6      def getQValue(self, state, action):
7          return self.values[state, action]
8
9      def computeValueFromQValues(self, state):
10         actions = self.getLegalActions(state)
```

```python
            if not actions:
                return 0
            return self.getQValue(state, self.getPolicy(state))

    def computeActionFromQValues(self, state):
        actions = self.getLegalActions(state)
        if not actions:
            return 0
        return max(actions, key=lambda action: self.getQValue(state, action))

    def getAction(self, state):
        legalActions = self.getLegalActions(state)
        action = None
        return action

    def update(self, state, action, nextState, reward):
        newValue = (1 - self.alpha) * self.getQValue(state, action) + self.alpha * (reward +
        self.values[state, action] = newValue

    def getPolicy(self, state):
        return self.computeActionFromQValues(state)

    def getValue(self, state):
        return self.computeValueFromQValues(state)
```