



**UNIVERSITATEA
TEHNICĂ
DIN CLUJ-NAPOCA**

Reinforcement learning

Inteligența Artificială

Autori: Adriana Mihnea

Grupa: 30236

FACULTATEA DE AUTOMATICA
SI CALCULATOARE

3 ianuarie 2024

Cuprins

1	Question 1: Value Iteration	2
2	Question 3: Policies	3
2.1	Prefer the close exit (+1), risking the cliff (-10)	3
2.2	Prefer the close exit (+1), but avoiding the cliff (-10)	3
2.3	Prefer the distant exit (+10), risking the cliff (-10)	4
2.4	Prefer the distant exit (+10), avoiding the cliff (-10)	4
2.5	Avoid both exits and the cliff (so an episode should never terminate)	4
3	Question 5: Q-Learning (Neimplementat)	4

1 Question 1: Value Iteration

```
1 def runValueIteration(self):
2     for _ in range(self.iterations):
3         newValues = util.Counter()
4         for state in self.mdp.getStates():
5             if not self.mdp.isTerminal(state):
6                 # Compute the Q-value for each action
7                 Q_values = [self.computeQValueFromValues(state, action) for
8                             ↪ action in self.mdp.getPossibleActions(state)]
9                 newValues[state] = max(Q_values)
10        self.values = newValues
```

RunValueIteration este metoda principala ce contine algoritmul pentru un anumit numar de iteratii (self.iterations). Initializeaza un nou set de valori (newValues), dictionar care va pastra valorile actualizate pentru fiecare stare. Pentru fiecare stare in MDP:

1. Daca starea nu este terminala:
 - calculeaza Q-value pentru fiecare actiune posibila din starea curenta folosind metoda computeQValueFromValues
 - Q-value reprezinta recompensa cumulativa asteptata in urma luarii unei anumite actiuni din starea curenta si urmand apoi politica optima
 - noua valoare pentru stare este setata ca fiind Q-value maxim din toate actiunile posibile
2. In momentul in care se completeaza toate iteratiile, valorile agentului vor fi completate cu noile valori calculate

```
1 def computeQValueFromValues(self, state, action):
2     qValue = sum(prob * (self.mdp.getReward(state, action, nextState) +
3     ↪ self.discount * self.values[nextState])
4     for nextState, prob in
5     ↪ self.mdp.getTransitionStatesAndProbs(state, action))
6     return qValue
```

ComputeQValueFromValues calculeaza Q-value pentru o pereche data de stare-actiune, folosind functia curenta de valoare (self.values). Itereaza toate posibilele stari urmatoare si probabilitatile lor de tranzitie (prob) date in starea si actiunea curenta. Pentru fiecare stare urmatoare, calculeaza recompensa imediata (self.mdp.getReward) si viitoarea valoare discounted (self.discount * self.values[nextState]). Q-value este suma acestor valori pentru toate starile urmatoare posibile.

```
1 def computeActionFromValues(self, state):
2     if self.mdp.isTerminal(state):
3         return None
4
5     possibleActions = self.mdp.getPossibleActions(state)
6     return max(possibleActions, key=lambda action:
7     ↪ self.computeQValueFromValues(state, action))
```

ComputeActionFromValues calculeaza cea mai buna actiune a unei stari date bazat pe functia de valoare curenta. Daca starea este terminala, returneaza "None" deoarece nu

exista actiuni legale in stari terminale. Altfel, ia lista actiunilor posibile pentru starea curenta si returneaza actiunea care maximizeaza Q-value pentru perechea stare-actiune data folosind metoda `computeQValueFromValues`. Argumentul "key" ii specifica functiei folosite sa extraga o cheie de comparare (in acest caz, Q-value) de la fiecare element.

2 Question 3: Policies

Acest exercitiu presupune alegerea urmatoarelor 3 valori pentru a alege cai optime in diferite conditii:

- **Discount** determina importanta recompenselor viitoare (mai aproape de 1) in comparatie cu recompensele imediate (mai aproape de 0).
- **Noise**: o valoare mai mare creste probabilitatea unor actiuni aleatorii, in timp ce o valoare mai mica pastreaza actiunile deterministice. Aceasta valoare este importanta cand vine vorba de a lua decizia de a urma sau nu cai riscante
- **Living Reward** reprezinta recompensa (sau penalizarea) pe care o primeste agentul pentru fiecare pas ce nu face parte dintr-o stare terminala. Aceasta valoare afecteaza tendinta agentului de a explora, a evita sau a prioritiza anumite cai

Mai jos sunt prezentate conditiile care trebuie indeplinite:

2.1 Prefer the close exit (+1), risking the cliff (-10)

```
1 def question3a():
2     answerDiscount = 0.9
3     answerNoise = 0.2
4     answerLivingReward = -1.9
5     return answerDiscount, answerNoise, answerLivingReward
```

Am setat discount-ul la 0.9 pentru a face agentul să acorde o importanță semnificativă recompenselor imediate, dar fără a neglija complet cele pe termen lung. Am setat noise-ul la 0.2 astfel incat agentul sa nu se indeparteze prea mult de la calea cea mai scurta. Am setat living reward-ul la -1.9 astfel incat sa se aleaga cu prioritate calea cea mai scurta, orice valoare mai mica de atat facand testul sa esueze.

2.2 Prefer the close exit (+1), but avoiding the cliff (-10)

```
1 def question3b():
2     answerDiscount = 0.5
3     answerNoise = 0.2
4     answerLivingReward = -1
5     return answerDiscount, answerNoise, answerLivingReward
```

Am setat discount-ul la 0.5 pentru a acorda atentie recompensei apropiate decât celei pe termen lung. Am setat noise la 0.2 pentru ca agentul sa faca actiuni aleatoare, dar nu suficient de multe incat sa ajunga pe o cale riscanta. Am setat living reward la -1 pentru a face agentul sa ajunga cat mai rapid la starea finala.

2.3 Prefer the distant exit (+10), risking the cliff (-10)

```
1 def question3c():
2     answerDiscount = 0.95
3     answerNoise = 0.05
4     answerLivingReward = -0.1
5     return answerDiscount, answerNoise, answerLivingReward
```

Am setat discount-ul la 0.95 pentru a acorda atentie recompensei indepartate. Am redus noise la 0.05 pentru a face acțiunile agentului mai deterministe și a reduce explorarea aleatorie. Am setat living reward la -0.1 pentru a oferi o mică penalizare pentru timpul petrecut în stări non-terminale, fără a încuraja explorarea inutilă.

2.4 Prefer the distant exit (+10), avoiding the cliff (-10)

```
1 def question3d():
2     answerDiscount = 0.8
3     answerNoise = 0.4
4     answerLivingReward = -1
5     return answerDiscount, answerNoise, answerLivingReward
```

Am setat discount-ul la 0.8 pentru a echilibra importanța recompenselor imediate și cele pe termen lung. Am crescut noise la 0.4 pentru a permite explorarea aleatorie. Totuși, am setat living reward la -1 pentru a încuraja ajungerea rapidă la stările terminale.

2.5 Avoid both exits and the cliff (so an episode should never terminate)

```
1 def question3e():
2     answerDiscount = 0.8
3     answerNoise = 0.4
4     answerLivingReward = 2
5     return answerDiscount, answerNoise, answerLivingReward
```

Am menținut discount-ul la 0.8 pentru a echilibra importanța recompenselor imediate și cele pe termen lung. Am menținut noise la 0.4 pentru a permite explorarea aleatorie. Am setat living reward la 2 pentru a oferi o recompensă mare pentru timpul petrecut în stări non-terminale, încurajând astfel explorarea extinsă.

3 Question 5: Q-Learning (Neimplementat)