



**UNIVERSITATEA
TEHNICĂ
DIN CLUJ-NAPOCA**

Reinforcement learning

Inteligența Artificială

Darius-Iulian Stan

Grupa: 30231

FACULTATEA DE AUTOMATICA
SI CALCULATOARE

18 Ianuarie 2023

Cuprins

1	Introducere	2
2	Q1: Value Iteration	2
3	Q2: Bridge Crossing Analysis	3
4	Q3: Policies	4
5	Q5: Q-Learning	5
6	Q6: Epsilon Greedy	6
7	Q7: Bridge Crossing Revisited	6
8	Q8: Q-Learning and Pacman	7

1 Introducere

În cadrul acestui proiect, în primă fază, implementăm Value iteration și analizăm parametrii modelului MDP (Markov Decision Process), el fiind caracterizat astfel:

- Un set de stări S
- Un set de acțiuni A (pentru fiecare stare din S)
- O funcție de tranziție $T(s,a,s')$, care reprezintă probabilitatea ca, din starea s , dorind să facem acțiunea a , ajungem în starea s'
- o funcție de recompense $R(s,a,s')$: recompense pentru că agentul încă este în viață (pot fi negative) și recompense de terminare a jocului (ori foarte mari, ori foarte mici, depinzând de caracterul stării de finish)
- O stare inițială și, eventual, una de exit

Ne aflăm într-un mediu non-deterministic, asta însemnând că dacă alegem să facem acțiunea a , nu este garantat că outcomu-ul acesteia este cel așteptat. (Spre exemplu, dacă vrem să mergem din poziția curentă în Nord, s-ar putea, de fapt, să ajungem în poziția din Est). Astfel, definim zgomotul ("noise") ca fiind probabilitatea ca agentul să facă o altă acțiune decât cea cerută.

Formula de bază a acestui proiect este cea a lui Bellman, care definește valoarea optimală a unei stări ca fiind maximul dintre valorile Q -stărilor ce pot fi atinse făcând orice acțiune posibilă din starea dată. Q -valoarea, la rândul ei, este definită, pentru o stare s și o acțiune a , calculând suma produsurilor dintre probabilitatea de a ajunge într-o anumită stare s' în urma efectuării acțiunii a în starea s , și suma recompensei acestei tranziții (s - a - s') cu valoarea stării respective* un discount anume.

2 Q1: Value Iteration

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

Pe baza formulei de recurență de mai sus, provenind din formula lui Bellman caracterizand funcția pentru valori optimale, am implementat metodele: `computeActionFromValues`, `computeQValueFromValues`, `runValueIteration` din clasa `ValueIterationAgent`.

cod:

```
1     def runValueIteration(self):
2
3         states = self.mdp.getStates()
4         for iteration in range(self.iterations):
5             current_values = util.Counter()
6             for s in states:
7                 if not self.mdp.isTerminal(s):
8                     current_values[s] = self.computeQValueFromValues(s, self.getAction(s))
9             self.values=current_values
10
11     def computeQValueFromValues(self, state, action):
12
13         transitions = self.mdp.getTransitionStatesAndProbs(state,action)
14         if transitions is None:
15             return 0
16         Qvalue = 0
```

```

17     for(s, T) in transitions:
18         R = self.mdp.getReward(state, action, s)
19         Qvalue += T*(R + self.discount*self.values[s])
20     return Qvalue
21
22 def computeActionFromValues(self, state):
23
24     if self.mdp.isTerminal(state):
25         return None
26     best = {"action": None, "value": None}
27     actions = self.mdp.getPossibleActions(state)
28     for action in actions:
29         qVal = self.computeQValueFromValues(state, action)
30         if best["action"] is None or best["value"] < qVal:
31             best = {"action": action, "value": qVal}
32     return best["action"]

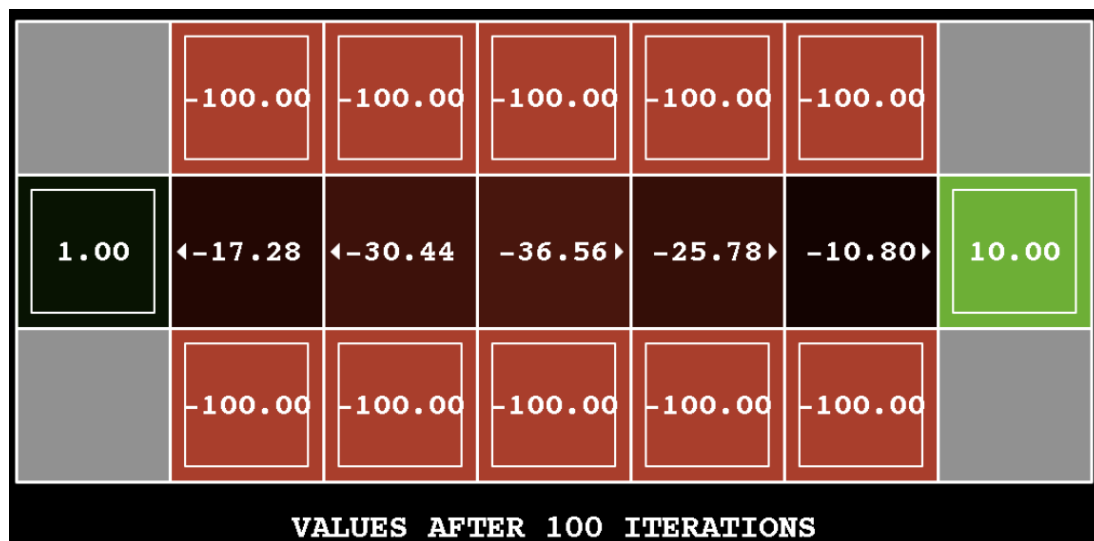
```

În `self.values` avem valorile pentru tabla noastră, ele inițial fiind zero (nu s-a efectuat nicio iterație, $k = 0$). Astfel, la fiecare iterație, pentru fiecare stare, dacă nu este terminală, actualizăm valoarea sa conform formulei, folosind datele de la iterația anterioară.

`computeActionFromValues` returnează cea mai bună acțiune pe care o pot face din starea dată (`state`), iar `computeQValueFromValues` crează suma din formula dată, pentru o anumită stare curentă și acțiune dată. Astfel, în `runValueIteration`, apelând `computeQValueFromValues` pentru starea curentă și acțiunea returnată de `computeActionFromValues`, nu mai trebuie să calculăm maximum dintre aceste sume, întrucât facem suma strict pentru cea mai bună acțiune.

3 Q2: Bridge Crossing Analysis

În cadrul acestui task, se dă următorul grid, reprezentând un "pod", trebuind să modificăm doar unul dintre parametrii dați (`noise` - probabilitatea să fac altă acțiune decât cea aleasă și `discount` - γ), ei fiind `noise = 0.2` și `discount = 0.9`.

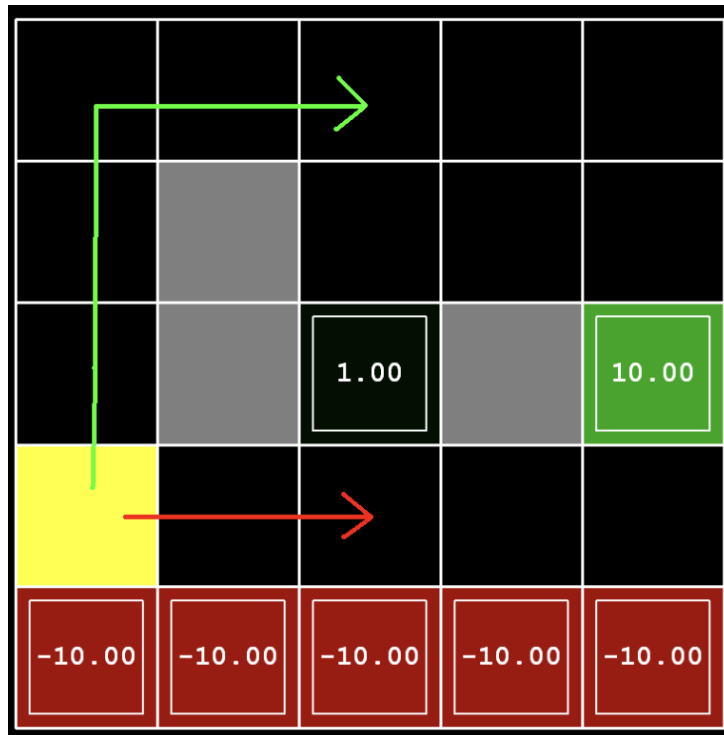


Conform formulei, pentru că, din orice stare componentă podului, pentru a-l traversa cu succes, dacă `noise`-ul ar fi 0.0, astfel că mediul ar fi unul deterministic, în care fiecare acțiune

duce la starea aleasă, evident că podul ar fi traversat cu succes (presupunând că discount-ul este totuși pozitiv, întrucât unul negativ ar inversa polaritatea valorilor optime în formula Q-valorilor, astfel că cele mai desirabile ar fi cele cu -100, podul ne mai fiind traversat). Totuși, din mai multe încercări, am observat că pentru discount-ul de 0.9, podul este traversat cu valoarea lui noise fiind dată din intervalul $[0.0, 0.016]$.

4 Q3: Policies

Asemănător cu task-ul anterior, dat fiind grid-ul alăturat, trebuie să dăm valori parametrilor (noise, discount, livingReward) astfel încât să avem urmoarele cazuri:



- Preferă finalState-ul 1.00, riscând să cazi în prăpastie: evident, pentru ca din patratul galben, să aleg ruta spre Est, în detrimentul celei din Nord, voi alege un noise extrem de mic, iar pentru a merge la cei mai apropiat finalState, setând livingReward-ul la o valoare negativă (-1, în cazul nostru), agentul nu va tenta să meargă până la finalState-ul 10.00, întrucât livingReward-ul este unul negativ, iar discount-ul, fiind și el mic, nu are cum să echivaleze (produsul dintre el și 10.00) valorile cumulate de livingRewards.
- Preferă finalState-ul 1.00, dar evitând prăpastia: pentru a alege ruta spre Nord, este suficient ca noise-ul să fie 0.3 (dacă este mai mare, atunci, cel mai probabil, ajungem într-o altă stare decât cea dorită, deci agentul va alege chiar acțiunea opusă, pentru a crește șansele să merg în direcția mea dorită inițial), livingReward-ul este pus pe -3, însă o valoare mai mare (în modul) poate duce la decizia agentului de a se risca să parcurgă calea către 1 de lângă prăpastie. Totodată, discount-ul de 0.4 ajută pentru echilibrarea livingRewardului, întrucât influențează agentul să meargă totuși spre un finalState pozitiv, garantând că scorul final va fi mai mare decât -10.00 (dacă s-ar arunca în prăpastie direct)
- Preferă FinalState-ul 10.00, riscând prăpastia: Asemănător primul scenariu, noise-ul trebuie să fie unul mic, iar discount-ul unul mare (0.9) astfel încât să fie preferată ruta care duce cel mai rapid spre finalState. Mai mult, livingResult-ul trebuie să fie îndeajuns de

mare încât să evite finalState-ul 1.00, însă nu prea mare pentru a eclipsa discount-ul și a alege ruta spre Nord.

- Preferă FinalState-ul 10, dar fără a risca prăpastia: LivingReward-ul unul mare, astfel încât agentul să dorească să viziteze cât mai multe stări, Noise-ul unul mediu (0.4), pentru a nu se dori parcurgerea căii de lângă prăpastie, iar discount-ul trebuie să fie destul de mare astfel încât să ajung, totuși, în finalState 20.0.
- : jocul nu se termină niciodată: evident, livingRewardul trebuie să eclipseze produsul dintre discounted și valorile finalState-urilor pozitive. Deci, livingState-ul-o valoare mare -, discount-ul unul mic, aproape de 0 (sau chiar 0). Noise-ul trebuie să fie, la fel, 0.4, pentru a nu se risca alegerea căii de lângă prăpastie, pentru că ar risca terminarea jocului (căderea în infern).

5 Q5: Q-Learning

Dat fiind că, în contextul subpunctelor anterioare, agentul nostru nu învață pe baza mediului în care se află în mod treptat, ci își crează harta acțiunilor corespunzătoare stărilor bazată pe valorile $V(s)$, lucru ce în lumea reală nu este posibil, pentru a crea totuși o strategie în care agentul să experimenteze fiecare stare și potențialele sale outcome-uri, implementăm Q-learning-ul. După cum îi spune și numele, el este bazat pe q-state-uri, nu pe valorile V ale stărilor propriu zise. Astfel, nu știm cu siguranță dacă o acțiune este întradevăr desirabilă sau nu, însă, prin efectuarea unui număr mare de iterații, cresc șansele de a primi toate variantele posibile pentru un Q-state, predominând, evident, cele cu probabilitatea mai mare. Astfel, pornind de la formula lui Bellman pentru Q-value, obținem:

$$Q(s, a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot \left(r + \gamma \cdot \max_{a'} Q(s', a') \right) \quad (1)$$

Pentru o vizibilitate mai sporită, vom nota suma din paranteză cu sample.

astfel, din formula de mai sus putem observa că, cu ajutorul lui alpha, care reprezintă rata de învățare, definim q-state-ul ca fiind media dintre valoarea pe care o avea deja din outcome-urile precedente și outcome-ul nou găsit (sample). Practic, Q-state-ul este de fapt media tuturor outcome-urilor găsite până în momentul de față cu noul outcome.

Astfel, în clasa `QLearningAgent`, am implementat metodele: `computeValueFromQValues` (returnează cea mai bună valoare a Q-state-urilor din starea s , pentru fiecare acțiune posibilă), `computeActionFromQValues` (returnează cea mai bună acțiune din state-ul dat), iar `update` actualizează Q-valoarea (state,action) conform formulei recursive de mai sus.

```
1 def getQValue(self, state, action):
2     return self.Qvalue[(state, action)]
3
4
5 def computeValueFromQValues(self, state):
6
7     actions = self.getLegalActions(state)
8     if not actions:
9         return 0
10    valueQAS = max([self.getQValue(state,action) for action in actions])
11    return valueQAS
12
```

```

13     def computeActionFromQValues(self, state):
14
15         actions = self.getLegalActions(state)
16         if not actions:
17             return 0
18         action = max(actions, key=lambda action: self.getQValue(state, action),
19                     default=None)
20         return action
21     def update(self, state, action, nextState, reward):
22
23         sample = reward + self.discount*self.getQValue(nextState,
24                 self.computeActionFromQValues(nextState))
25         self.Qvalue[(state, action)] = (1 - self.alpha)*self.getQValue(state, action) +
26                 self.alpha*(sample)

```

6 Q6: Epsilon Greedy

Pentru ca lumea în care antrenăm agentul să fie cât mai imprevizibilă, construim metoda `getAction()` în următorul fel: pentru starea în care ne aflăm, alegem o acțiune în mod aleatoriu în epsilon la sută din cazuri, în celelalte cazuri, selectăm cea mai bună acțiune disponibilă din starea respectivă, alegere făcută pe cunoștințele dobândite până în prezent.

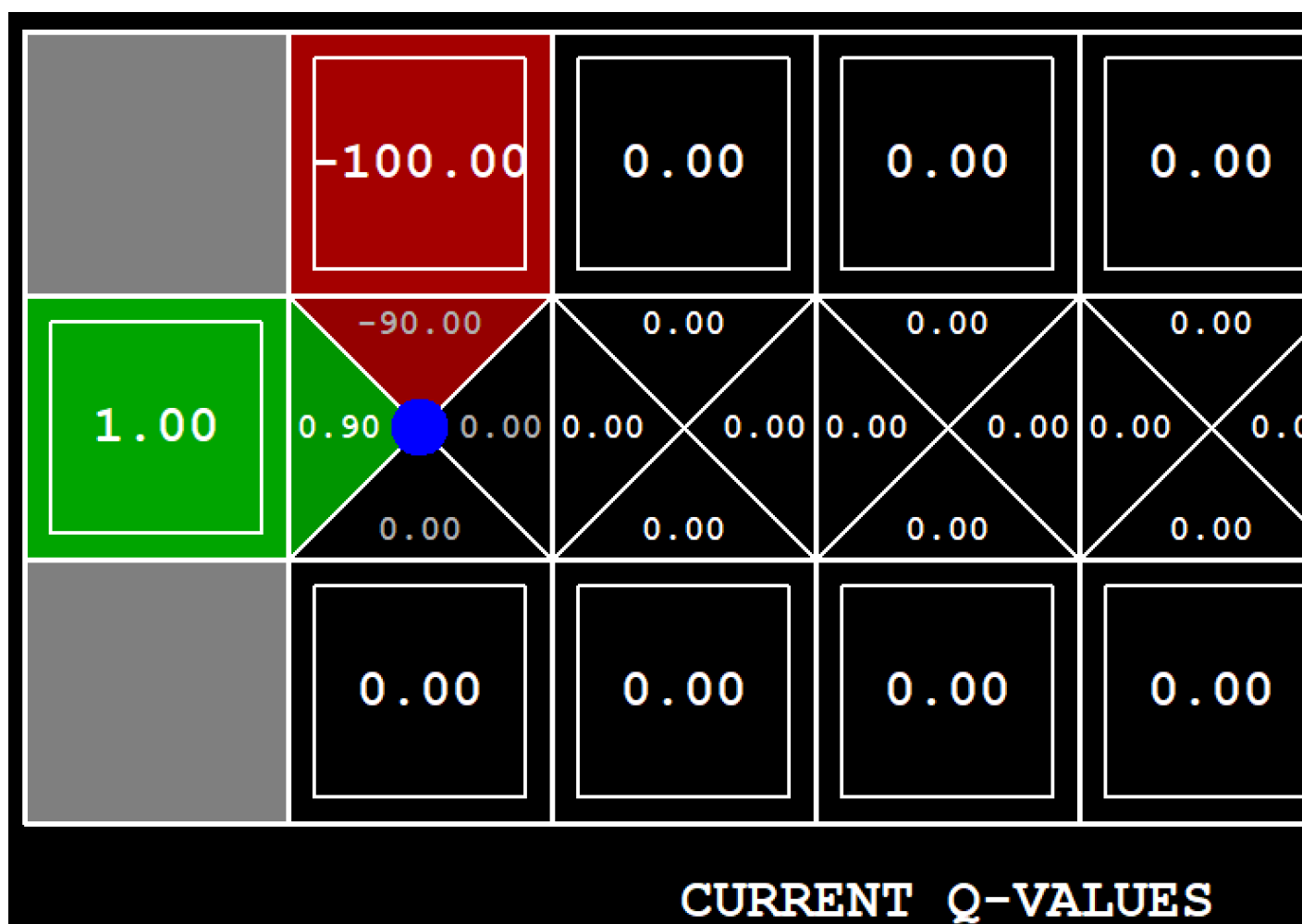
```

1     def getAction(self, state):
2         # Pick Action
3         legalActions = self.getLegalActions(state)
4         action = None
5
6         if not legalActions:
7             return None
8         if util.flipCoin(self.epsilon):
9             action = random.choice(legalActions)
10        else:
11            action = self.getPolicy(state)
12
13        return action

```

7 Q7: Bridge Crossing Revisited

Revenind la problema podului din Q2, vrem să verificăm dacă, pentru Q-learning, există valori pentru epsilon și alpha (learning rate) astfel încât un anumit grid (cel din imagine), să convergă după 50 de k iterații. Însă, răspunsul este negativ! Dat fiind că epsilon de câte ori fac o acțiune neprevăzută, afectează direct proporțional spațiul stărilor care sunt vizitate (de exemplu, pentru epsilon 0, la un moment dat va exista bucla continuă de Est - West dintre casutele). Dacă epsilon este 1, atunci, într-adevăr, spațiul stărilor exploreate se extinde, însă nu îndeajuns astfel încât să se acopere, cu valorile corespunzătoare, fiecare celulă.



8 Q8: Q-Learning and Pacman

Cu ajutorul Q-learn-ului, Pacman reușește, după 2000 de teste pentru învățare a mediului, să aplice cu succes politica lui bazată de experiență, evitând fantomele și luând recompensa.