

# **Învățarea de către un agent a jocului Snake folosind algoritmul Q-Learning**

**Universitatea Tehnică „Gheorghe Asachi” din Iași**



**Facultatea de Automatică și Calculatoare din Iași**



## **Autori**

Daniel Rotariu

Eduard Pușcașu

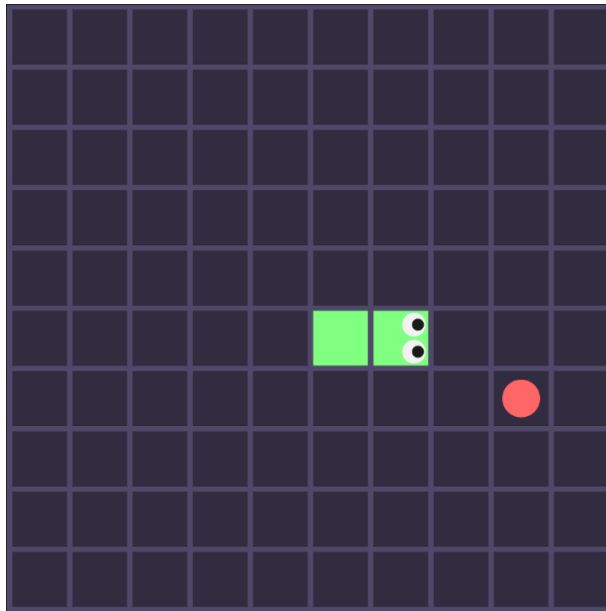
## **1. Introducere**

Snake este un joc în care jucătorul manevrează capătul unei linii în creștere, adesea cu tematica unui șarpe. Jucătorul trebuie să împiedice șarpele să se ciocnească atât de alte obstacole, cât și de el însuși. Scorul unui jucător este dat de numărul de bucăți de mâncare consumate. Fiecare mâncare consumată face ca șarpele să crească în lungime, făcând astfel mai dificilă manevrarea acestuia.

În această lucrare ne propunem să folosim algoritmul Q-Learning pentru a antrena un agent la jocul Snake, scopul acestui agent fiind de a obține un scor cât mai mare.

## **2. Descrierea problemei**

Avem un mediu 2D format din 3 componente principale: șarpele, mâncarea și marginile spațiului în care șarpele se poate mișca.

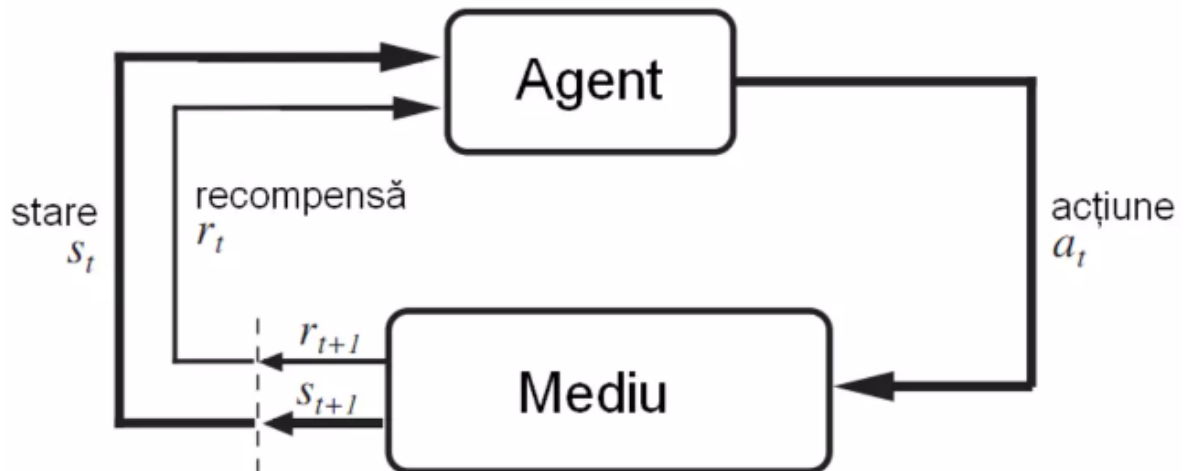


În această problemă, scopul nostru este de a antrena agentul astfel încât acesta să facă deciziile corecte în ceea ce privește mișcarea șarpelui pentru a se îndrepta spre bucata de mâncare existentă la un moment dat în timp și să evite să iasă în afara spațiului permis sau să se ciocnească de el însuși.

## **3. Aspecte teoretice privind algoritmul**

- **Învățare cu întărire**

Este o arie a învățării automate care are ca obiectiv maximizarea recompensei primite de către un agent care ia o serie de decizii într-un mediu dinamic.



Pașii făcuți de un agent sunt următorii:

- Observa mediul
- Folosește o strategie pentru a decide ce acțiune va alege în funcție de starea mediului
- În urma acțiunii alese anterior ajunge într-o nouă stare
- Primește o recompensa sau o penalizare
- Învăță din experiența sa și își modifica strategia
- Itereaza pana cand găsește o strategie optimă

Exista doua tipuri de algoritmi de învățare prin întărire:

- Bazati pe model: învață modelul de tranziții și recompense și îl folosește pentru a descoperi politica optimă
- Fara model: descoperă politica optimă fără a învăța sau estima modelul

### • Algoritmul Q-Learning

Algoritmul Q-Learning este un algoritm de învățare prin întărire fără model. Acest algoritm va găsi cea mai buna secvență de acțiuni pornind de la starea în care se afla agentul. În funcție de starea mediului în care se afla agentul, acesta va alege următoarea acțiune fără a se baza pe o politică existentă, ci bazându-se doar pe experiența sa.

După un număr de  $n$  pași în viitor, agentul va fi nevoit să aleaga următoarea sa acțiune, ponderea recompensei acestei acțiuni va fi  $\gamma^n$ , unde  $\gamma$  (discount factor) reprezintă un număr în intervalul  $[0,1]$  și are ca efect alegerea unei acțiuni care va aduce o recompensă mai mare într-un număr cât mai mic de pași.

Așadar, algoritmul are o funcție care calculează recompensa primită de agent în funcție de starea în care se afla și de acțiunea pe care o poate lua.

$$Q : S \times A \rightarrow R$$

Valorile acestei funcții sunt stocate într-o structură de date numită Q-Table.

Înainte ca procesul de învățare să înceapă, această tabelă este inițializată cu valori arbitrare fixe. Apoi, la orice moment de timp  $t$ , agentul care se află în starea  $s_t$  alege o acțiune  $a_t$ , observă recompensa  $r_t$  primită și ajunge într-o nouă stare  $s_{t+1}$  care poate depinde de starea anterioară și de acțiunea aleasă. În final tabelă Q este actualizată. Ecuația care stă la baza acestui algoritm este *Ecuația lui Bellman*.

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

unde  $r_t$  reprezintă recompensa primită de către agent trecând din starea  $s_t$  în starea  $s_{t+1}$  prin acțiunea  $a_t$ .

Un episod al algoritmului se termină când starea  $s_{t+1}$  este o stare terminală. Pentru toate stările finale,  $s_f$  tabelă Q nu se actualizează.

- **Influența parametrilor**

- Learning rate -  $\alpha$ : determină în ce măsură informațiile noi învățate suprascriu informațiile vechi. Valoarea 0 pentru acest parametru va determina ca agentul să nu învețe nimic, bazându-se doar pe cunoștințele proprii, iar valoarea 1 va determina agentul să folosească doar informațiile cele mai recente.
- Discount factor -  $\gamma$ : determină importanța viitoarelor recompense. Valoarea 0 va determina agentul ca acesta să aleagă mereu recompensele curente și să le ignore pe cele din stările viitoare, iar valoarea 1 va determina agentul să ia în considerare cea mai mare recompensă cumulată din stările viitoare.

#### **4. Modalitatea de rezolvare**

Performanța acestui algoritm este puternic influențată de modul în care ne definim stările mediului. Un număr prea mic de stări nu ar cuprinde toate situațiile în care s-ar putea regăsi agentul, asta ar duce la alegerea unei acțiuni nepotrivite și la obținerea în medie a unui scor mic. Un număr prea mare de stări poate duce la un consum inutil de resurse.

Abordarea aleasă de noi definește o stare ca o combinație a următoarelor caracteristici:

- Direcția în care se deplasează șarpele: aceasta poate lua 4 valori (stanga, dreapta, sus, jos)
- Direcția mănării raportată la poziția capului șarpelui: inițial, valorile posibile au fost stanga, dreapta, sus, jos. Acest lucru a dus la un comportament greșit al agentului atunci când mănarea nu se afla în linie dreaptă cu capul șarpelui. Spre

exemplu, cand mancarea se afla in stanga sus fata de acesta, agentul trebuie sa aleaga dacă direcția mănării este în stanga sau în sus și uneori va alege acțiuni ce îl vor îndepărta de aceasta. Prin urmare, am ales ca valorile posibile pentru aceasta caracteristica sa fie stanga, dreapta, sus, jos, stanga-sus, stanga-jos, dreapta-sus, dreapta-jos.

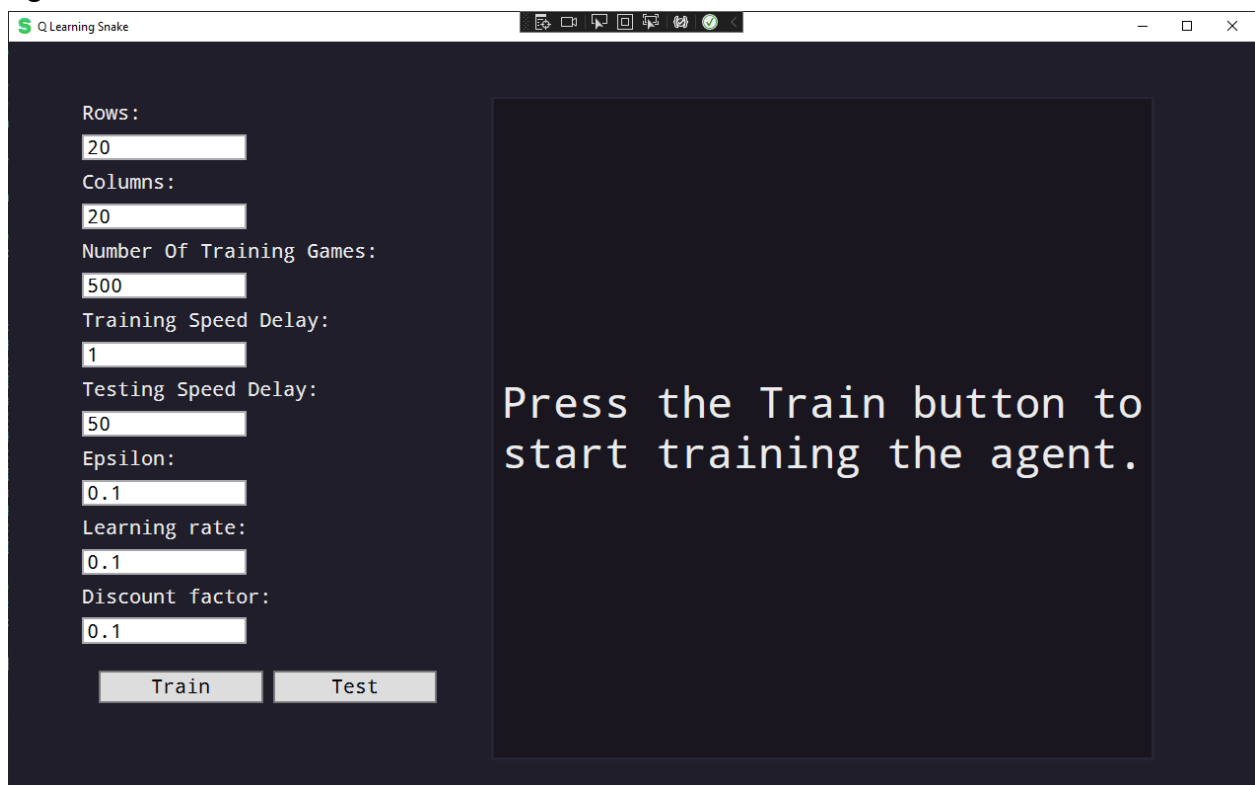
- Pericol în direcțiile stanga, dreapta, sus, jos: deoarece putem avea mai multe pericole simultan, vom avea 4 variabile ce pot lua valorile adevărat sau fals.

În total avem  $4*8*2*2*2*2 = 512$  stări posibile. Pentru fiecare stare în care se afla, agentul poate alege între 3 acțiuni: să își păstreze direcția, să schimbe direcția și să meargă în sus sau în jos (dacă direcția curentă de deplasare este stanga sau dreapta) sau să schimbe direcția și să meargă în stanga sau în dreapta (dacă direcția curentă de deplasare este în sus sau în jos).

Recompensele stabilite pentru fiecare rezultat al acțiunilor alese de agent sunt următoarele:

- Dacă se apropie de mancare: +1
- Dacă se îndepărtează de mancare: -1
- Dacă ajunge la mancare: +10
- Dacă se lovește de un obstacol și pierde jocul: -1000

Din interfața putem selecta dimensiunea spațiului în care se poate mișca șarpele, numărul de jocuri pentru care dorim să antrenăm agentul, viteza de antrenare/testare, și parametrii algoritmului.



La apăsarea butonului Train, se inițializează tabela Q. Am ales un dicționar pentru a stoca valorile în care cheia o reprezintă combinația caracteristicilor descrise anterior (adică starea), iar valoarea este o listă în care indicii reprezintă acțiunile și valorile reprezintă recompensele asociate acțiunilor respective.

```
foreach (Direction headDirection in headDirections)
{
    foreach (Direction foodDirection in foodDirections)
    {
        foreach (bool dangerLeft in dangerLeftOptions)
        {
            foreach (bool dangerRight in dangerRightOptions)
            {
                foreach (bool dangerUp in dangerUpOptions)
                {
                    foreach (bool dangerDown in dangerDownOptions)
                    {
                        String key = $"{headDirection.GetHashCode()}, {foodDirection.GetHashCode()}, {dangerLeft}, {dangerRight}, {dangerUp}, {dangerDown}";
                        qTable.Add(key, new double[3] { 0, 0, 0 });
                    }
                }
            }
        }
    }
}
```

La fiecare pas, calculăm starea în care se afla agentul, iar la alegerea unei acțiuni din acea stare, dacă agentul se antrenează vom folosi tehnica *Epsilon-Greedy*. În majoritatea cazurilor, agentul va alege cele mai bune acțiuni în funcție de experiența sa, acest lucru poate duce la ignorarea unei acțiuni neexplorate care poate aduce o recompensă mai mare. Deci folosim un parametru  $\epsilon$  ca agentul să aleagă în  $\epsilon\%$  cazuri o acțiune aleatoare.

```
↑ reference
public Direction GetNextMove(GameState gameState, bool isTraining)
{
    int action;
    double prob = random.NextDouble();
    String currentState = GetState(gameState);

    if (isTraining && prob <= epsilon)
    {
        action = GetRandomAction();
    }
    else
    {
        action = GetBestAction(currentState);
    }

    Direction direction = GetActionDirection(gameState, action);
    GameState nextGameState = gameState.DeepClone();

    nextGameState.ChangeDirection(direction);
    nextGameState.Move();

    UpdateQTable(gameState, action, nextGameState);

    return direction;
}
```

După ce agentul alege acțiunea, vom calcula noua stare în care acesta ajunge și vom actualiza tabela Q folosind *Ecuația lui Bellman* și recompensele corespunzătoare.

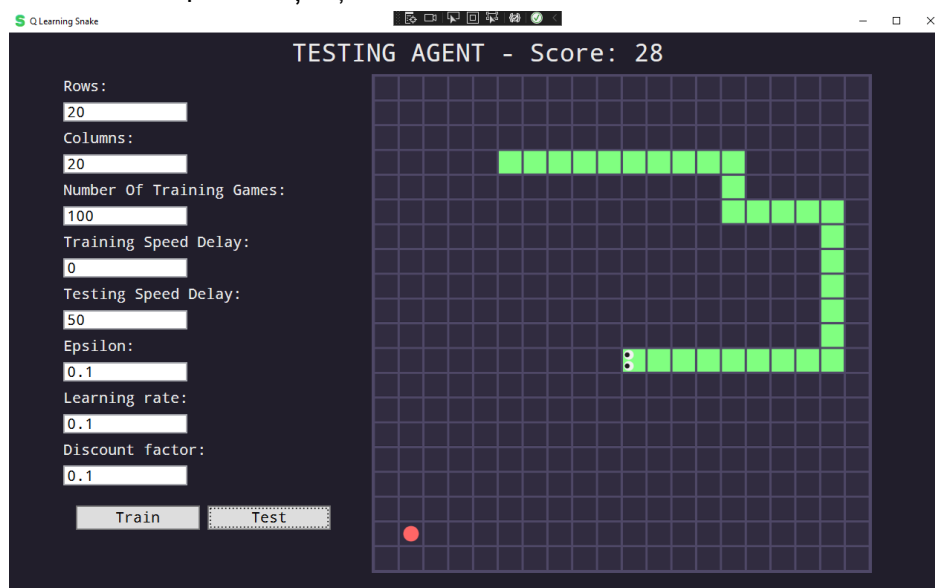
```
private void UpdateQTable(GameState gameState, int action, GameState nextGameState)
{
    String oldState = GetState(gameState);
    String newState = GetState(nextGameState);
    int reward;

    if (nextGameState.GameOver)
    {
        reward = GAME_OVER_REWARD;

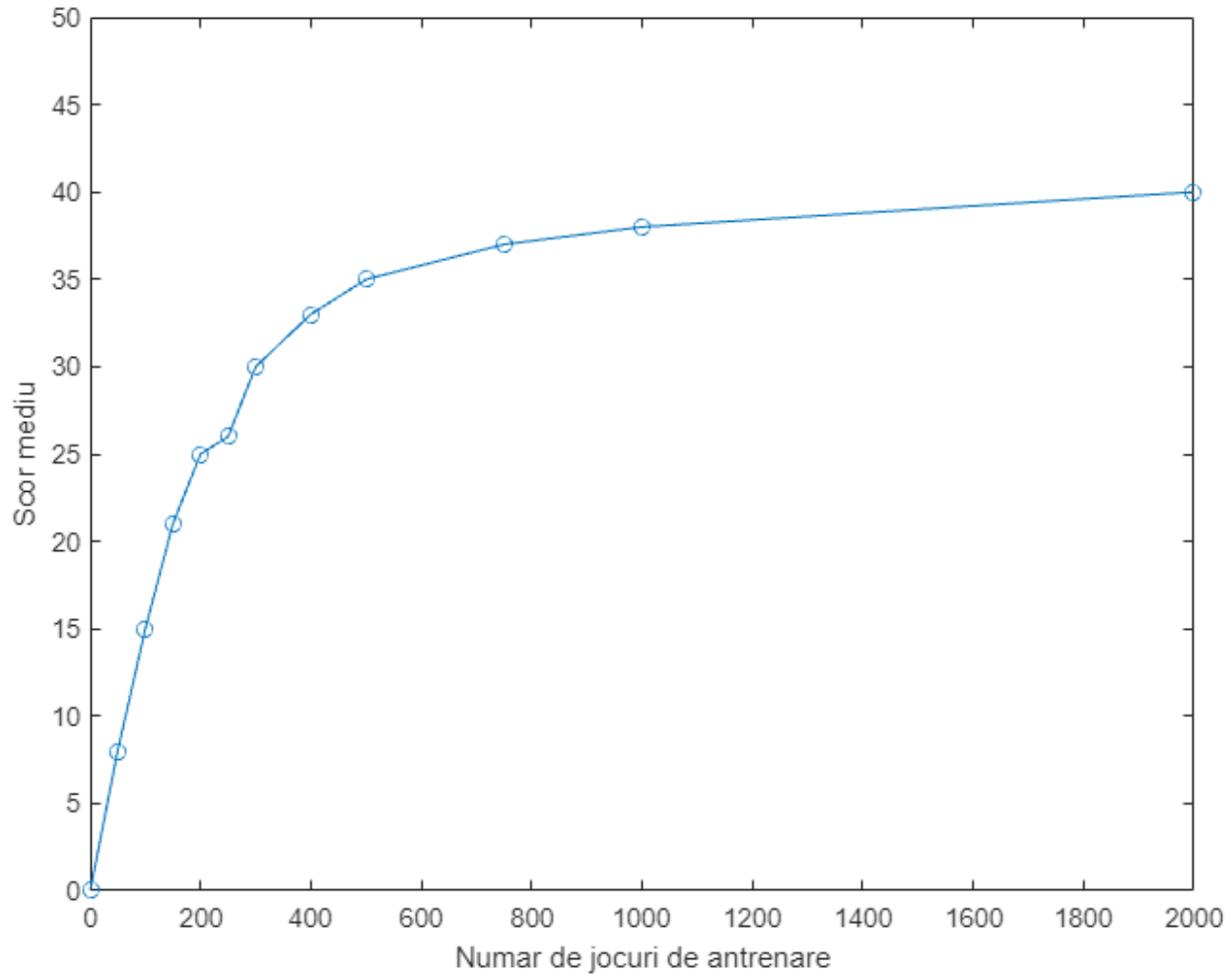
        // Bellman Equation - No Future State
        qTable[oldState][action] = qTable[oldState][action] + learningRate * (reward - qTable[oldState][action]);
    }
    else
    {
        if (nextGameState.HeadPosition().Equals(gameState.foodPosition))
        {
            reward = EAT_FOOD_REWARD;
        }
        else if (nextGameState.RowDistanceFromFood() < gameState.RowDistanceFromFood() ||
            nextGameState.ColDistanceFromFood() < gameState.ColDistanceFromFood())
        {
            reward = MOVE_CLOSER_TO_FOOD_REWARD;
        }
        else
        {
            reward = MOVE_AWAY_FROM_FOOD_REWARD;
        }

        // Bellman Equation
        int nextStateBestAction = GetBestAction(newState);
        qTable[oldState][action] = qTable[oldState][action] + learningRate *
            (reward + discountFactor * qTable[newState][nextStateBestAction] - qTable[oldState][action]);
    }
}
```

La finalul perioadei de antrenare, folosind butonul Test putem testa agentul, în aceasta faza acesta se va baza numai pe cunoștințele dobândite în faza de antrenare.



## 5. Rezultate obtinute

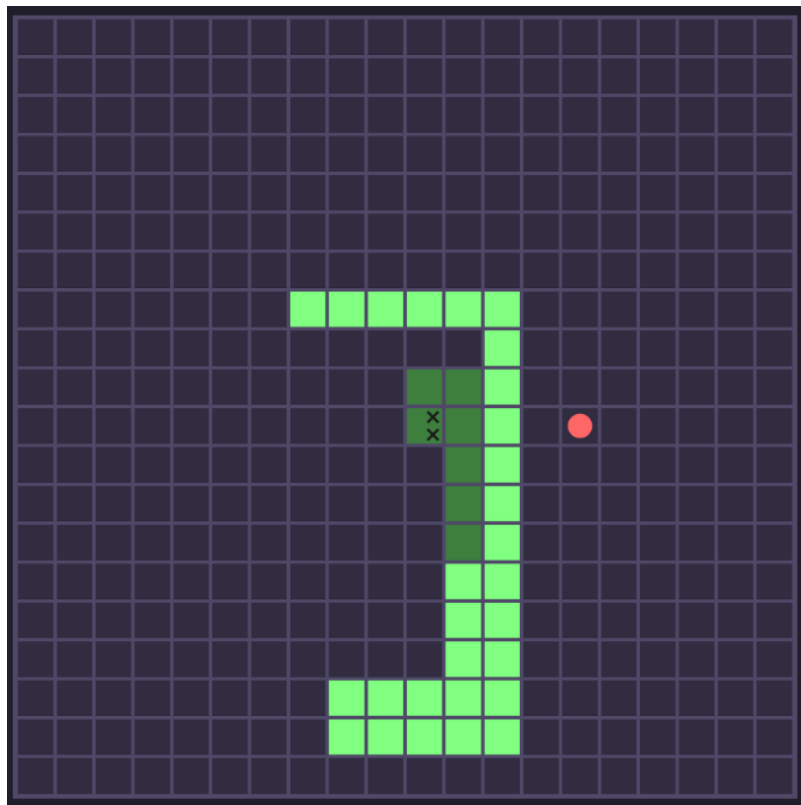


Observam ca scorul mediu pe care reușește să îl obțină agentul nostru crește proporțional cu numărul de jocuri de antrenare și atinge un platou după aproximativ 500-600 de jocuri. După acest prag e nevoie de un număr de jocuri mai mare de 10.000 (timpul necesar antrenării fiind de peste 3 ore) pentru a sesiza o creștere considerabilă a scorului.

## 6. Concluzii

În implementarea propusă de noi, agentul nu va juca niciodată un joc perfect. Acesta reușește să se apropie de mâncare și să evite pericolele imediate, însă nu este suficient de inteligent pentru a evita să se blocheze în propriul corp.





În scenariul de mai sus, agentul ar trebui sa aleaga direcția stanga ca următoarea mișcare a sarpelui, pentru a-si evita corpul, însă alege sa nu se indeparteze de mancare, blocandu-se astfel. Ar trebui modificat spațiul stărilor astfel incat sa includem și acest caz pentru ca agentul sa poată învăța să îl evite.

## 7. Contributii

Implementare joc Snake	Eduard Pușcașu
Implementare interfata	Daniel Rotariu
Stabilire stari, recompense, parametri	Daniel Rotariu, Eduard Pușcașu
Implementare algoritm Q-Learning	Daniel Rotariu
Testare agent in diverse situatii	Eduard Pușcașu
Documentatie - Sectiunile 1, 2, 3	Daniel Rotariu
Documentatie - Sectiunile 4, 5, 6	Eduard Pușcașu

## 8. **Bibliografie**

- <https://en.wikipedia.org/wiki/Q-learning>
- <https://www.slideshare.net/FlorinLeon/invatarea-cu-intarire>
- <https://towardsdatascience.com/q-learning-algorithm-from-explanation-to-implementation-cdbeda2ea187>
- <https://www.cse.unsw.edu.au/~cs9417ml/RL1/algorithms.html>
- <http://cs229.stanford.edu/proj2016spr/report/060.pdf>
- <https://openreview.net/pdf?id=iu2XOJ45cxo>
- <https://towardsdatascience.com/teaching-a-computer-how-to-play-snake-with-q-learning-93d0a316ddc0>