

# Rapport de projet - Apprentissage par renforcement

## Participants :

- Abdallah Abdelsadek
- Adam Tawfik

## Études de cas : Environnements Simple et Tic Tac Toe

Ce rapport décrit les choix d'implémentation et les approches utilisées pour résoudre deux études de cas dans le cadre du projet **MN(AA)** : un environnement simplifié (**Simple**) et le jeu du Morpion (**TicTacToe**) à l'aide d'algorithmes d'apprentissage par renforcement.

### 1. Environnement Simple

Pour cet environnement statique à un seul état, nous avons testé trois algorithmes de bandits multi-bras :

- Epsilon-Greedy
- Gradient Bandit
- UCB (Upper Confidence Bound)

Chaque algorithme a été intégré dans une architecture d'agent et a permis d'apprendre les actions optimales à long terme en maximisant la récompense moyenne. L'agent n'interagit pas avec un adversaire, ce qui permet une convergence rapide.

### 2. Environnement Tic Tac Toe

Pour le jeu du Morpion, nous avons utilisé les algorithmes suivants :

- Q-Learning

---

#### Algorithm 5 Q-Learning

---

**Require:** Taux d'apprentissage  $\alpha > 0$ , facteur d'actualisation  $\gamma \in (0, 1)$

**Ensure:** Fonction d'action-valeur  $Q(s, a)$  convergeant vers  $Q^*(s, a)$

1: **Initialiser**  $Q(s, a)$  arbitrairement pour tous les états  $s$  et actions  $a$

2: **for** chaque épisode **do**

3:   **Observer** ou **choisir** un état initial  $s$

4:   **while**  $s$  n'est pas terminal **do**

5:     **Choisir** une action  $a$  en fonction de  $Q$  (avec une exploration  $\varepsilon$ -greedy)

6:     **Exécuter** l'action  $a$  et **observer** la récompense  $r$  et l'état suivant  $s'$

7:     **Mettre à jour**  $Q(s, a)$ :

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

8:     **Affecter**  $s \leftarrow s'$

9:   **end while**

10: **end for**

11: **Retourner** la table  $Q(s, a)$

---

- Value Iteration

---

**Algorithm 4** Value Iteration pour estimer  $\pi^*$ 


---

**Require:** Seuil de convergence  $\theta > 0$ , facteur d'actualisation  $\gamma$   
**Ensure:** Fonction de valeur  $V(s)$  pour tout  $s \in S$  et politique optimale  $\pi^*$

- 1: Initialiser  $V(s)$  arbitrairement pour tout  $s \in S^+$ , avec  $V(\text{terminal}) \leftarrow 0$
- 2: **repeat**
- 3:      $\Delta \leftarrow 0$
- 4:     **for** chaque état  $s \in S$  **do**
- 5:          $v \leftarrow V(s)$
- 6:          $V(s) \leftarrow \max_a \sum_{s',r} p(s',r | s,a) [r + \gamma V(s')]$
- 7:          $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- 8:     **end for**
- 9: **until**  $\Delta < \theta$
- 10: **Retourner** la politique déterministe  $\pi^*$  telle que, pour tout  $s \in S$ :

$$\pi^*(s) = \arg \max_a \sum_{s',r} p(s',r | s,a) [r + \gamma V(s')].$$

- Policy Iteration

---

**Algorithm 3** Iterative Policy Evaluation with Policy improvement

---

**Require:**  $S, A, p(s',r | s,a), \gamma \in [0,1], \theta > 0$   
**Ensure:** Politique optimale  $\pi^*$  et fonction de valeur optimale  $V^*$

- 1: Initialiser arbitrairement  $V(s)$  et  $\pi(s)$  pour tout  $s \in S$ , avec  $V(\text{terminal}) \leftarrow 0$
- 2: **repeat**
- 3:     **repeat**
- 4:          $\Delta \leftarrow 0$
- 5:         **for** chaque état  $s \in S$  **do**
- 6:              $v \leftarrow V(s)$
- 7:              $V(s) \leftarrow \sum_a \pi(a | s) \sum_{s',r} p(s',r | s,a) [r + \gamma V(s')]$
- 8:              $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- 9:         **end for**
- 10:     **until**  $\Delta < \theta$
- 11:     policy-stable  $\leftarrow$  true
- 12:     **for** chaque état  $s \in S$  **do**
- 13:          $a_{\text{old}} \leftarrow \pi(s)$
- 14:          $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r | s,a) [r + \gamma V(s')]$
- 15:         **if**  $a_{\text{old}} \neq \pi(s)$  **then**
- 16:             policy-stable  $\leftarrow$  false
- 17:         **end if**
- 18:     **end for**
- 19: **until** policy-stable est vraie
- 20: **return**  $\pi^*$  et  $V^*$

---

Nous avons défini une fonction de récompense spécifique :

**+1:** quand le joueur 0 gagne, **-1:** quand le joueur 1 gagne, **0:** en cas de match nul, et **-0.25:** pour tout autre coup. Cela permet d'encourager des victoires rapides et décourager les actions non optimales.

Le joueur 0 (**agent**) apprend seul dans un environnement partiellement coopératif. À la fin de l'apprentissage, il est capable de battre des adversaires aléatoires ou coopératifs. Nous avons intégré un mode de jeu [ **Agent vs Human** ] et [ **Agent vs Agent** ].

## Conclusion

Les environnements et les algorithmes ont été choisis pour illustrer des cas simples et complexes. Le système développé est modulaire, ce qui facilite l'ajout d'autres environnements et algorithmes. Nos choix de fonction de récompense et de stratégie d'adversaire ont fortement influencé la performance de l'apprentissage.