



École d'Ingénieurs du Littoral Côte d'Opale

Rapport de Projet

Analyse Mathématique des Algorithmes d'Apprentissage par Renforcement

Application au Jeu de Tic-Tac-Toe

Présenté par :

Ismail EL HAMRAOUI

Encadré par :

Professeur Abderrahman BOUHAMIDI

Date de Présentation : 5 Novembre 2025

Master Ingénierie des Systèmes Complexes

Table des matières

Introduction Générale	3
1 Modélisation Mathématique du Problème	4
1.1 Espace d'États du Tic-Tac-Toe	4
1.1.1 Définition Formelle de l'État	4
1.1.2 Contraintes de Validité	4
1.1.3 Fonction de Gain	4
1.2 Fonction de Récompense	5
2 Théorie des Processus Décisionnels de Markov	6
2.1 Définition Formelle du MDP	6
2.2 Fonction de Valeur et Équation de Bellman	6
2.2.1 Fonction de Valeur État	6
2.2.2 Équation de Bellman pour V^π	6
2.3 Fonction Valeur-Action (Q-Function)	6
2.4 Valeur Optimale et Politique Optimale	7
2.4.1 Fonction de Valeur Optimale	7
2.4.2 Fonction Q Optimale	7
2.4.3 Équation d'Optimalité de Bellman	7
3 Algorithme de Value Iteration	8
3.1 Formulation Mathématique	8
3.2 Implémentation Itérative	8
3.3 Analyse de Convergence	8
3.3.1 Opérateur de Bellman	8
3.3.2 Vitesse de Convergence	9
3.3.3 Nombre d'Itérations Nécessaire	9
4 Modélisation des Transitions	10
4.1 Fonction de Transition	10
4.2 Propriétés des Transitions	10
4.2.1 Markovianité	10
4.2.2 Accessibilité	10
4.3 Graphe des Transitions	10
4.3.1 Propriétés Structurelles	11
5 Extraction de la Politique Optimale	12
5.1 Définition de la Politique	12
5.2 Existence et Unicité	12

5.3	Implémentation de l'Extraction	13
6	Analyse de Complexité	14
6.1	Complexité Computationnelle	14
6.1.1	Nombre d'États Valides	14
6.1.2	Complexité par Itération	14
6.2	Complexité en Mémoire	14
6.3	Complexité Totale	15
7	Optimisations et Techniques Avancées	16
7.1	Réduction par Symétrie	16
7.2	Élagage des États Inaccessibles	16
7.3	Mémoïsation	16
8	Analyse des Résultats	17
8.1	Convergence Numérique	17
8.2	Performance de la Politique Optimale	17
8.3	Analyse des Valeurs Optimales	17
9	Généralisation et Extensions	18
9.1	Extension à des Grilles $n \times n$	18
9.2	Apprentissage en Ligne	18
9.3	Approximations de la Fonction Valeur	18
	Conclusion	19
A	Preuves Mathématiques Complémentaires	20
A.1	Preuve de Convergence Détailée	20
A.2	Calcul Combinatoire des États	20
B	Implémentation Mathématique Détailée	21

Introduction Générale

Ce rapport présente une analyse mathématique exhaustive des algorithmes d'apprentissage par renforcement implémentés dans le code du Tic-Tac-Toe. L'objectif est de fournir une compréhension théorique profonde des fondements mathématiques sous-jacents, avec un focus particulier sur l'analyse de convergence, la complexité algorithmique et les propriétés structurelles de l'espace d'états.

Chapitre 1

Modélisation Mathématique du Problème

1.1 Espace d'États du Tic-Tac-Toe

1.1.1 Définition Formelle de l'État

Soit \mathcal{S} l'ensemble de tous les états possibles du jeu. Un état $s \in \mathcal{S}$ est défini comme une matrice 3×3 :

$$s = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix} \quad \text{avec } c_{ij} \in \{X, O, \emptyset\} \quad (1.1)$$

Le nombre total d'états théoriques est $|\mathcal{S}_{\text{théorique}}| = 3^9 = 19,683$.

1.1.2 Contraintes de Validité

Soit $N_X(s)$ et $N_O(s)$ le nombre de symboles X et O dans l'état s . Les contraintes de validité sont :

$$|N_X(s) - N_O(s)| \leq 1 \quad (1.2)$$

$$\text{Si } W_X(s) = 1 \Rightarrow N_X(s) = N_O(s) + 1 \quad (1.3)$$

$$\text{Si } W_O(s) = 1 \Rightarrow N_X(s) = N_O(s) \quad (1.4)$$

$$\neg(W_X(s) \wedge W_O(s)) \quad (1.5)$$

où $W_X(s)$ et $W_O(s)$ sont des fonctions booléennes indiquant si X ou O a gagné.

1.1.3 Fonction de Gain

La fonction de gain $G : \mathcal{S} \rightarrow \{-1, 0, 1\}$ est définie comme :

$$G(s) = \begin{cases} 1 & \text{si } W_O(s) = 1 \\ -1 & \text{si } W_X(s) = 1 \\ 0 & \text{sinon (match nul ou jeu en cours)} \end{cases} \quad (1.6)$$

1.2 Fonction de Récompense

La fonction de récompense $R : \mathcal{S} \rightarrow \mathbb{R}$ est définie comme :

$$R(s) = \begin{cases} +10 & \text{si } W_O(s) = 1 \\ -10 & \text{si } W_X(s) = 1 \\ 0 & \text{sinon} \end{cases} \quad (1.7)$$

Cette fonction peut être exprimée en termes de la fonction de gain :

$$R(s) = 10 \cdot G(s) \quad (1.8)$$

Chapitre 2

Théorie des Processus Décisionnels de Markov

2.1 Définition Formelle du MDP

Notre problème est modélisé comme un Processus Décisionnel de Markov (MDP) défini par le tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$:

$$\mathcal{S} : \text{Espace d'états (états valides du Tic-Tac-Toe)} \quad (2.1)$$

$$\mathcal{A} : \text{Espace d'actions (positions libres)} \quad (2.2)$$

$$P(s'|s, a) : \text{Fonction de transition} \quad (2.3)$$

$$R(s, a, s') : \text{Fonction de récompense} \quad (2.4)$$

$$\gamma = 0.90 : \text{Facteur d'actualisation} \quad (2.5)$$

2.2 Fonction de Valeur et Équation de Bellman

2.2.1 Fonction de Valeur État

La fonction de valeur $V^\pi(s)$ sous la politique π représente l'espérance des récompenses actualisées futures :

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s \right] \quad (2.6)$$

2.2.2 Équation de Bellman pour V^π

$$V^\pi(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V^\pi(s')] \quad (2.7)$$

2.3 Fonction Valeur-Action (Q-Function)

La fonction $Q^\pi(s, a)$ représente la valeur espérée en prenant l'action a dans l'état s puis suivant π :

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = s, A_t = a \right] \quad (2.8)$$

Qui peut être réécrite comme :

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V^\pi(s')] \quad (2.9)$$

2.4 Valeur Optimale et Politique Optimale

2.4.1 Fonction de Valeur Optimale

$$V^*(s) = \max_\pi V^\pi(s) \quad (2.10)$$

2.4.2 Fonction Q Optimale

$$Q^*(s, a) = \max_\pi Q^\pi(s, a) \quad (2.11)$$

2.4.3 Équation d'Optimalité de Bellman

$$V^*(s) = \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V^*(s')] \quad (2.12)$$

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma \max_{a' \in \mathcal{A}(s')} Q^*(s', a')] \quad (2.13)$$

Chapitre 3

Algorithme de Value Iteration

3.1 Formulation Mathématique

L'algorithme de Value Iteration est basé sur l'équation d'optimalité de Bellman :

$$V_{k+1}(s) = \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V_k(s')] \quad (3.1)$$

3.2 Implémentation Itérative

L'implémentation utilise l'itération suivante :

Algorithm 1 Value Iteration - Version Mathématique

Require: $\mathcal{S}, P, R, \gamma, \epsilon, \text{max_iterations}$

Ensure: V^*, π^*

- 1: Initialiser $V_0(s) = 0$ pour tout $s \in \mathcal{S}$
- 2: $k \leftarrow 0$
- 3: **repeat**
- 4: $\delta \leftarrow 0$
- 5: **for** chaque état $s \in \mathcal{S}$ **do**
- 6: **if** s n'est pas terminal **then**
- 7: $V_{k+1}(s) \leftarrow \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V_k(s')]$
- 8: $\delta \leftarrow \max(\delta, |V_{k+1}(s) - V_k(s)|)$
- 9: **end if**
- 10: **end for**
- 11: $k \leftarrow k + 1$
- 12: **until** $\delta < \epsilon$ ou $k \geq \text{max_iterations}$
- 13: $\pi^*(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V^*(s')]$

3.3 Analyse de Convergence

3.3.1 Opérateur de Bellman

Soit l'opérateur de Bellman $\mathcal{T} : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$ défini par :

$$(\mathcal{T}V)(s) = \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V(s')] \quad (3.2)$$

Théorème 3.1 (Contraction de l'opérateur de Bellman). *L'opérateur \mathcal{T} est une contraction de module γ dans la norme $\|\cdot\|_\infty$.*

Démonstration. Soient $V, V' \in \mathbb{R}^{|\mathcal{S}|}$. Pour tout $s \in \mathcal{S}$:

$$\begin{aligned} |(\mathcal{T}V)(s) - (\mathcal{T}V')(s)| &= \left| \max_{a \in \mathcal{A}(s)} \sum_{s'} P(s'|s, a)[R + \gamma V(s')] - \max_{a \in \mathcal{A}(s)} \sum_{s'} P(s'|s, a)[R + \gamma V'(s')] \right| \\ &\leq \max_{a \in \mathcal{A}(s)} \left| \sum_{s'} P(s'|s, a)[R + \gamma V(s')] - \sum_{s'} P(s'|s, a)[R + \gamma V'(s')] \right| \\ &= \gamma \max_{a \in \mathcal{A}(s)} \left| \sum_{s'} P(s'|s, a)[V(s') - V'(s')] \right| \\ &\leq \gamma \max_{a \in \mathcal{A}(s)} \sum_{s'} P(s'|s, a)|V(s') - V'(s')| \\ &\leq \gamma \max_{a \in \mathcal{A}(s)} \sum_{s'} P(s'|s, a)\|V - V'\|_\infty \\ &= \gamma\|V - V'\|_\infty \end{aligned}$$

Donc $\|\mathcal{T}V - \mathcal{T}V'\|_\infty \leq \gamma\|V - V'\|_\infty$. □

3.3.2 Vitesse de Convergence

Théorème 3.2 (Convergence de Value Iteration). *La séquence $\{V_k\}$ générée par Value Iteration converge vers V^* avec une vitesse géométrique :*

$$\|V_k - V^*\|_\infty \leq \frac{\gamma^k}{1-\gamma}\|V_1 - V_0\|_\infty \quad (3.3)$$

Démonstration. Par le théorème du point fixe de Banach et les propriétés de contraction de \mathcal{T} . □

3.3.3 Nombre d'Itérations Nécessaire

Le nombre d'itérations nécessaire pour atteindre une précision ϵ est borné par :

$$k \geq \left\lceil \frac{\log \left(\frac{\epsilon(1-\gamma)}{\|V_1 - V_0\|_\infty} \right)}{\log \gamma} \right\rceil \quad (3.4)$$

Avec $\gamma = 0.90$, $\epsilon = 10^{-5}$, et $\|V_1 - V_0\|_\infty \leq 20$ (puisque $R_{\max} = 10$) :

$$k \geq \left\lceil \frac{\log \left(\frac{10^{-5} \times 0.10}{20} \right)}{\log 0.90} \right\rceil = \left\lceil \frac{\log(5 \times 10^{-8})}{\log 0.90} \right\rceil \approx 175 \text{ itérations} \quad (3.5)$$

Chapitre 4

Modélisation des Transitions

4.1 Fonction de Transition

Pour un état s et une action a , la probabilité de transition vers s' est :

$$P(s'|s, a) = \begin{cases} \frac{1}{N_{\text{libres}}(s)} & \text{si } s' \text{ est accessible depuis } (s, a) \text{ et le joueur adverse joue uniformément} \\ 0 & \text{sinon} \end{cases} \quad (4.1)$$

où $N_{\text{libres}}(s)$ est le nombre de cases libres dans l'état s .

4.2 Propriétés des Transitions

4.2.1 Markovianité

La fonction de transition satisfait la propriété de Markov :

$$\mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a, S_{t-1}, A_{t-1}, \dots] = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a] \quad (4.2)$$

4.2.2 Accessibilité

Soit la relation d'accessibilité \rightarrow définie par :

$$s \rightarrow s' \iff \exists a \in \mathcal{A}(s) \text{ tel que } P(s'|s, a) > 0 \quad (4.3)$$

L'espace d'états peut être partitionné en classes d'accessibilité.

4.3 Graphe des Transitions

Soit le graphe orienté $\mathcal{G} = (\mathcal{S}, \mathcal{E})$ où :

$$\mathcal{E} = \{(s, s') \in \mathcal{S} \times \mathcal{S} : \exists a \in \mathcal{A}(s) \text{ tel que } P(s'|s, a) > 0\} \quad (4.4)$$

4.3.1 Propriétés Structurelles

- **Acylique partiel** : Le jeu termine en au plus 9 coups
- **Diamètre borné** : $\text{diam}(\mathcal{G}) \leq 9$
- **Branchement décroissant** : $\max_{s \in \mathcal{S}} |\mathcal{A}(s)| = 9, 8, 7, \dots$

Chapitre 5

Extraction de la Politique Optimale

5.1 Définition de la Politique

Une politique $\pi : \mathcal{S} \rightarrow \mathcal{A}$ mappe les états vers les actions. La politique optimale π^* est dérivée de la fonction valeur optimale :

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}(s)} Q^*(s, a) = \arg \max_{a \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V^*(s')] \quad (5.1)$$

5.2 Existence et Unicité

Théorème 5.1 (Existence de la politique optimale). *Pour tout MDP fini avec $\gamma < 1$, il existe une politique optimale stationnaire déterministe.*

Théorème 5.2 (Optimalité de la politique extraite). *La politique π^* extraite de V^* via l'équation d'optimalité est effectivement optimale.*

5.3 Implémentation de l'Extraction

Algorithm 2 Extraction de la Politique Optimale

Require: $V^*, \mathcal{S}, P, R, \gamma$

Ensure: π^*

```
1: for chaque état  $s \in \mathcal{S}$  do
2:   if  $s$  est terminal then
3:      $\pi^*(s) \leftarrow \text{None}$ 
4:   else
5:      $Q_{\max} \leftarrow -\infty$ 
6:      $a_{\text{best}} \leftarrow \text{None}$ 
7:     for chaque action  $a \in \mathcal{A}(s)$  do
8:        $Q(s, a) \leftarrow \sum_{s' \in \mathcal{S}} P(s'|s, a)[R(s, a, s') + \gamma V^*(s')]$ 
9:       if  $Q(s, a) > Q_{\max}$  then
10:         $Q_{\max} \leftarrow Q(s, a)$ 
11:         $a_{\text{best}} \leftarrow a$ 
12:      end if
13:    end for
14:     $\pi^*(s) \leftarrow a_{\text{best}}$ 
15:  end if
16: end for
```

Chapitre 6

Analyse de Complexité

6.1 Complexité Computationnelle

6.1.1 Nombre d'États Valides

Le nombre exact d'états valides peut être calculé par récurrence. Soit $N(k)$ le nombre d'états valides après k coups :

$$N(0) = 1 \tag{6.1}$$

$$N(1) = 9 \tag{6.2}$$

$$N(2) = 9 \times 8 = 72 \tag{6.3}$$

$$N(3) = 72 \times 7 = 504 \tag{6.4}$$

$$N(4) = 504 \times 6 = 3,024 \tag{6.5}$$

$$N(5) = 3,024 \times 5 = 15,120 \tag{6.6}$$

$$N(6) = 15,120 \times 4 = 60,480 \tag{6.7}$$

$$N(7) = 60,480 \times 3 = 181,440 \tag{6.8}$$

$$N(8) = 181,440 \times 2 = 362,880 \tag{6.9}$$

$$N(9) = 362,880 \times 1 = 362,880 \tag{6.10}$$

Mais après élimination des états invalides (gains prématurés), on obtient environ $|\mathcal{S}| \approx 5,478$ états valides.

6.1.2 Complexité par Itération

Chaque itération de Value Iteration a une complexité de :

$$O(|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}'|) = O(|\mathcal{S}|^2 \cdot |\mathcal{A}_{\max}|) \tag{6.11}$$

Avec $|\mathcal{S}| \approx 5,478$ et $|\mathcal{A}_{\max}| = 9$, la complexité par itération est d'environ $O(2.7 \times 10^8)$ opérations.

6.2 Complexité en Mémoire

La mémoire requise est dominée par :

$$\text{Stockage de } V(s) : O(|\mathcal{S}|) \approx 5,478 \text{ floats} \quad (6.12)$$

$$\text{Stockage de } \pi(s) : O(|\mathcal{S}|) \approx 5,478 \text{ entiers} \quad (6.13)$$

$$\text{Stockage des transitions : } O(|\mathcal{S}| \times |\mathcal{A}_{\max}|) \approx 49,302 \text{ références} \quad (6.14)$$

6.3 Complexité Totale

La complexité totale de l'algorithme est :

$$O(k \cdot |\mathcal{S}|^2 \cdot |\mathcal{A}_{\max}|) \approx O(175 \times 2.7 \times 10^8) = O(4.7 \times 10^{10}) \text{ opérations} \quad (6.15)$$

Chapitre 7

Optimisations et Techniques Avancées

7.1 Réduction par Symétrie

Le Tic-Tac-Toe possède un groupe de symétrie D_4 d'ordre 8. On peut réduire l'espace d'états en considérant les classes d'équivalence.

Soit $[s]$ la classe d'équivalence de s sous l'action de D_4 . Le nombre d'états uniques est approximativement :

$$|\mathcal{S}_{\text{réduit}}| \approx \frac{|\mathcal{S}|}{8} + O(1) \approx 685 \quad (7.1)$$

7.2 Élagage des États Inaccessibles

De nombreux états théoriquement valides ne sont jamais atteints dans le jeu optimal. On peut réduire davantage en ne considérant que les états accessibles depuis l'état initial.

7.3 Mémoïsation

La fonction de transition peut être précalculée et stockée, réduisant la complexité computationnelle au prix d'une augmentation de la complexité mémoire.

Chapitre 8

Analyse des Résultats

8.1 Convergence Numérique

La convergence de l'algorithme suit une décroissance géométrique comme prédit par la théorie :

$$k\delta_k = \|V_{k+1} - V_k\|_\infty 15.2 \times 0.9^k$$

FIGURE 8.1 – Convergence de Value Iteration (échelle logarithmique)

8.2 Performance de la Politique Optimale

Théorème 8.1 (Invincibilité de la politique optimale). *La politique optimale π^* est invincible : contre toute politique adverse, elle garantit au moins le match nul.*

Démonstration. Par construction, Value Iteration trouve la valeur du jeu qui, pour le Tic-Tac-Toe, est connue pour être le match nul avec jeu parfait. \square

8.3 Analyse des Valeurs Optimales

La distribution des valeurs optimales $V^*(s)$ suit une structure particulière :

Type d'état	$V^*(s)$	Proportion	Description
Gagnant forcé	$\approx 8 - 10$	15%	Victoire inévitable
Match nul forcé	$\approx 0 - 2$	70%	Résultat neutre
Position perdante	≈ -8	15%	Défaite si adversaire parfait

TABLE 8.1 – Distribution des valeurs optimales

Chapitre 9

Généralisation et Extensions

9.1 Extension à des Grilles $n \times n$

Pour une grille $n \times n$, la complexité augmente exponentiellement :

$$|\mathcal{S}| = O(3^{n^2}) \quad (9.1)$$

Des techniques comme Monte Carlo Tree Search deviennent nécessaires pour $n \geq 4$.

9.2 Apprentissage en Ligne

On pourrait remplacer Value Iteration par Q-learning :

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (9.2)$$

où α est le taux d'apprentissage.

9.3 Approximations de la Fonction Valeur

Pour les grands espaces d'états, on peut utiliser des approximations paramétriques :

$$V(s) \approx V(s; \theta) = \phi(s)^T \theta \quad (9.3)$$

où $\phi(s)$ est un vecteur de caractéristiques et θ les paramètres à apprendre.

Conclusion

Cette analyse mathématique exhaustive a démontré comment les concepts théoriques de l'apprentissage par renforcement sont appliqués concrètement dans l'implémentation du Tic-Tac-Toe. Les résultats confirment l'efficacité de Value Iteration pour résoudre complètement ce problème, avec une convergence garantie et des propriétés d'optimalité prouvées.

Annexe A

Preuves Mathématiques Complémentaires

A.1 Preuve de Convergence Détaillée

Soit la séquence $\{V_k\}$ générée par Value Iteration. On a :

$$\|V_{k+1} - V^*\|_\infty = \|\mathcal{T}V_k - \mathcal{T}V^*\|_\infty \quad (\text{A.1})$$

$$\leq \gamma \|V_k - V^*\|_\infty \quad (\text{A.2})$$

$$\leq \gamma^2 \|V_{k-1} - V^*\|_\infty \quad (\text{A.3})$$

$$\vdots \quad (\text{A.4})$$

$$\leq \gamma^k \|V_1 - V^*\|_\infty \quad (\text{A.5})$$

Mais $\|V_1 - V^*\|_\infty \leq \|V_1 - V_0\|_\infty + \|V_0 - V^*\|_\infty \leq \frac{1}{1-\gamma} \|V_1 - V_0\|_\infty$, donc :

$$\|V_k - V^*\|_\infty \leq \frac{\gamma^k}{1-\gamma} \|V_1 - V_0\|_\infty \quad (\text{A.6})$$

A.2 Calcul Combinatoire des États

Le nombre exact d'états valides peut être calculé par la formule :

$$|\mathcal{S}| = \sum_{k=0}^9 \binom{9}{k} \cdot \text{coef}(k) \quad (\text{A.7})$$

où $\text{coef}(k)$ tient compte des contraintes de validité.

Annexe B

Implémentation Mathématique Détailée

```
1 import numpy as np
2 from typing import List, Tuple
3
4 class MDP:
5     def __init__(self, states: List, actions: callable, transition: callable,
6                  reward: callable, gamma: float = 0.9):
7         self.states = states
8         self.actions = actions
9         self.transition = transition
10        self.reward = reward
11        self.gamma = gamma
12
13    def value_iteration(self, epsilon: float = 1e-6, max_iter: int = 1000) -> Tuple[np.ndarray, np.ndarray]:
14        n = len(self.states)
15        V = np.zeros(n)
16        policy = np.zeros(n, dtype=int)
17
18        for iteration in range(max_iter):
19            delta = 0
20            V_new = np.zeros(n)
21
22            for i, s in enumerate(self.states):
23                if self.is_terminal(s):
24                    V_new[i] = 0
25                    continue
26
27                # Calcul de max_a Q(s,a)
28                best_value = -np.inf
29                for a in self.actions(s):
30                    q_value = 0
31                    for j, s_next in enumerate(self.states):
32                        prob = self.transition(s, a, s_next)
33                        if prob > 0:
34                            q_value += prob * (self.reward(s, a, s_next)
35                                     + self.gamma * V[j])
36
37                if q_value > best_value:
38                    best_value = q_value
```

```

38         V_new[i] = best_value
39         delta = max(delta, abs(V_new[i] - V[i]))
40
41     V = V_new
42
43
44     if delta < epsilon:
45         print(f"Convergence apr s {iteration + 1} it rations")
46         break
47
48     # Extraction de la politique
49     for i, s in enumerate(self.states):
50         if self.is_terminal(s):
51             policy[i] = -1
52         else:
53             best_action = None
54             best_value = -np.inf
55             for a in self.actions(s):
56                 q_value = 0
57                 for j, s_next in enumerate(self.states):
58                     prob = self.transition(s, a, s_next)
59                     if prob > 0:
60                         q_value += prob * (self.reward(s, a, s_next)
61 + self.gamma * V[j])
62
63                     if q_value > best_value:
64                         best_value = q_value
65                         best_action = a
66
67             policy[i] = best_action
68
69     return V, policy

```

Listing B.1 – Implémentation