

Apprentissage par renforcement (Reinforcement Learning (RL)) Approche : Temporal Difference (TD)

T. AL-ANI
A²SI-ESIEE-PARIS



A²SI

Plan de la présentation

- I. Introduction
- II. L'apprentissage par renforcement
- III. Temporal Difference

I. Introduction

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference
- VI. Conclusion et Perspectives

- Contexte du Projet
- Objectif de notre projet
 - Étude de la théorie de l'apprentissage par renforcement
 - Implémentation des méthodes de Temporal Difference

II. L'apprentissage par renforcement

I.	Introduction
II.	L'apprentissage par renforcement
III.	Programmation dynamique
IV.	Monte Carlo
V.	Temporal Difference
VI.	Conclusion et Perspectives

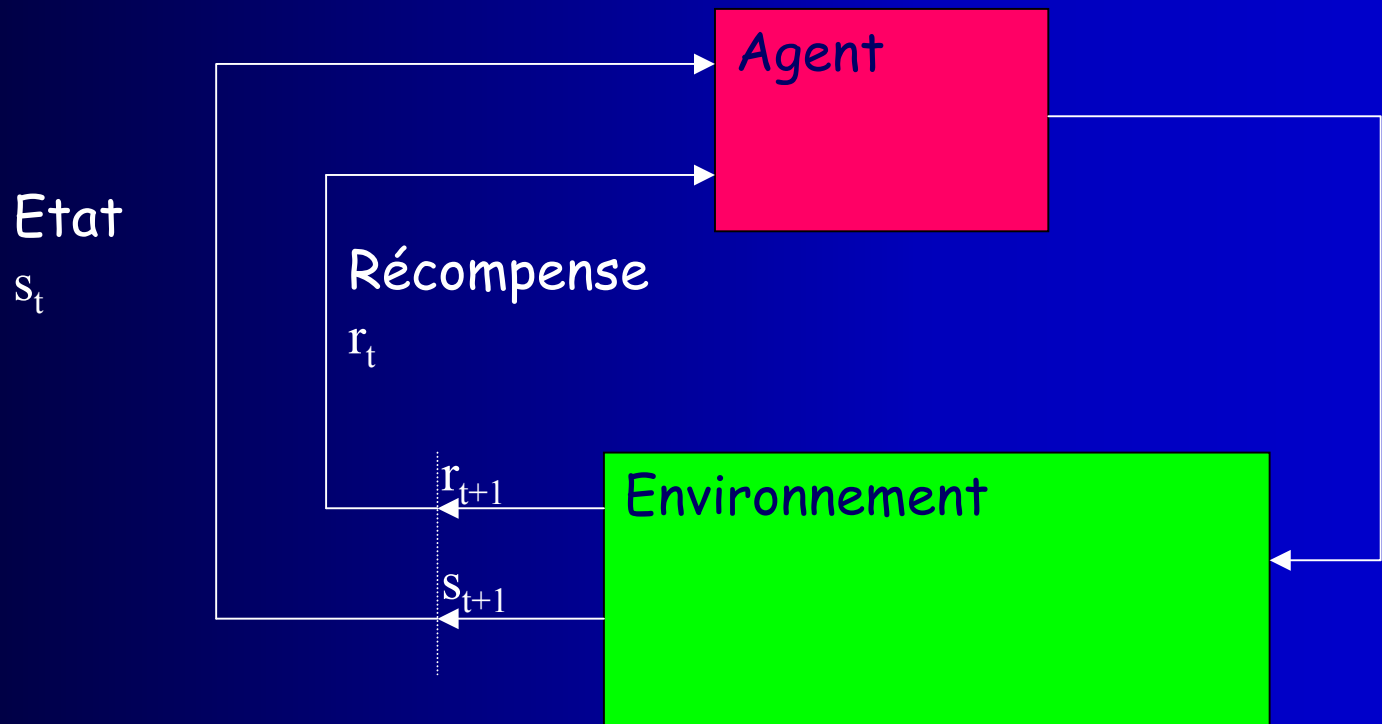
- Généralités

- Apprentissage par ‘essais / erreurs’
- Récompenses suivant les actions
- Équilibre entre Exploration et Exploitation

II. L'apprentissage par renforcement

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference
- VI. Conclusion et Perspectives

• Interaction Agent – Environnement



II. L'apprentissage par renforcement

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference
- VI. Conclusion et Perspectives

- Théorie

- Évaluation de politique

- Amélioration de politique

V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference**
- VI. Conclusion et Perspectives

– Aucun modèle de l'environnement est nécessaire, l'expérience suffit.

– L'apprentissage se fait avant de connaître la récompense finale voire sans la connaître.

– Les méthodes convergent.

V. Temporal Difference

I.	Introduction
II.	L'apprentissage par renforcement
III.	Programmation dynamique
IV.	Monte Carlo
V.	Temporal Difference
VI.	Conclusion et Perspectives

• Théorie: TD prédiction

1. Initialize $V(s)$ arbitrarily, π to the policy to be evaluated :

2. Repeat (for each episode) :

 Initialize s

 Repeat (for each step of episode) :

$a \leftarrow$ action given by π for s

 Take action a ; observe reward, r , and next state, s'

$V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$

$s \leftarrow s'$

 until s is terminal

ALGO. 5.1: Algorithme $TD(0)$ d'évaluation d'une politique

V. Temporal Difference

I.	Introduction
II.	L'apprentissage par renforcement
III.	Programmation dynamique
IV.	Monte Carlo
V.	Temporal Difference
VI.	Conclusion et Perspectives

• Théorie: TD Sarsa

```

1. Initialize  $Q(s, a)$  arbitrarily :

2. Repeat (for each episode) :
    Initialize  $s$ 
    Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat (for each step of episode) :
        Take action  $a$ , observe  $r, s'$ 
        Choose  $a'$  from  $s'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$ 
         $s \leftarrow s'; a \leftarrow a'$ 
    until  $s$  is terminal
    
```

ALGO. 5.2: Algorithme *Sarsa*

V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. **Temporal Difference**
- VI. Conclusion et Perspectives

• Théorie: TD Q learning

```

1. Initialize  $Q(s, a)$  arbitrarily :

2. Repeat (for each episode) :
    Initialize  $s$ 
    Repeat (for each step of episode) :
        Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
        Take action  $a$ , observe  $r, s'$ 
         $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ 
         $s \leftarrow s'$  ;
    until  $s$  is terminal
    
```

ALGO. 5.3: Algorithme *Q-learning*

V. Temporal Difference

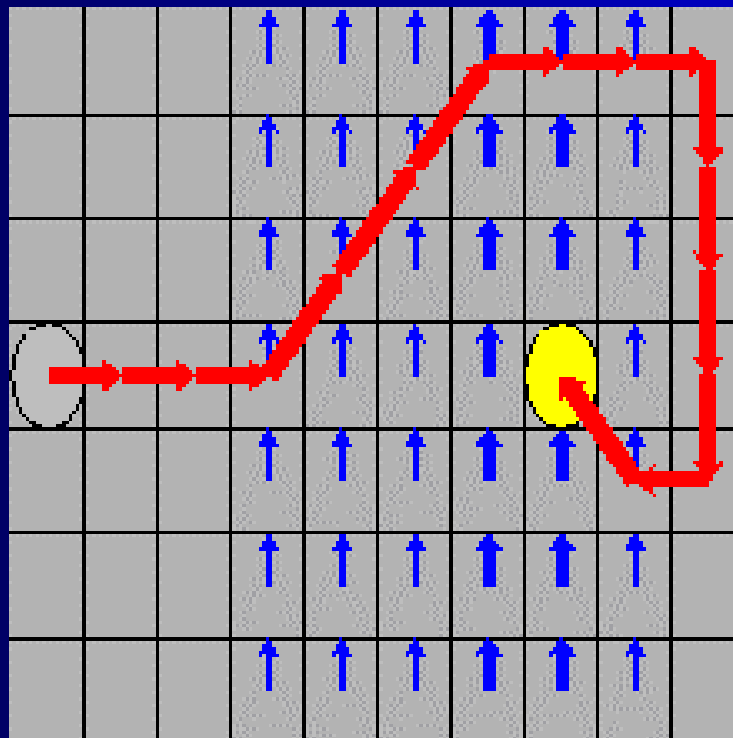
- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference**
- VI. Conclusion et Perspectives

- Exemple du Random Walk

V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. **Temporal Difference**
- VI. Conclusion et Perspectives

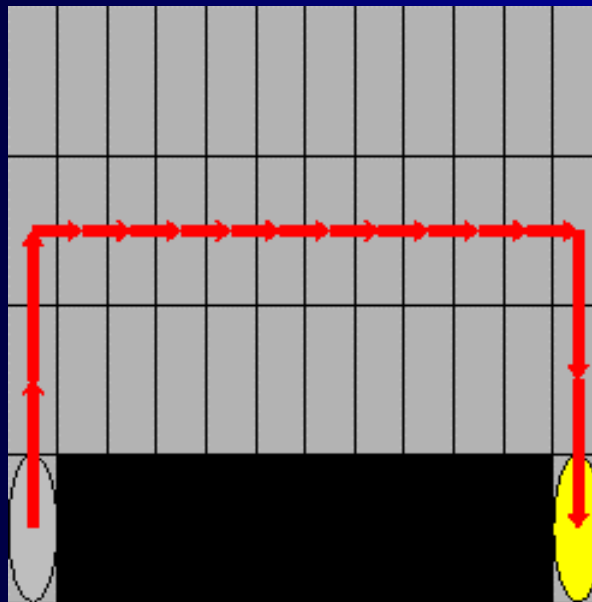
- Exemple 1 du *Windy Grid*
– Exemple du livre de Sutton



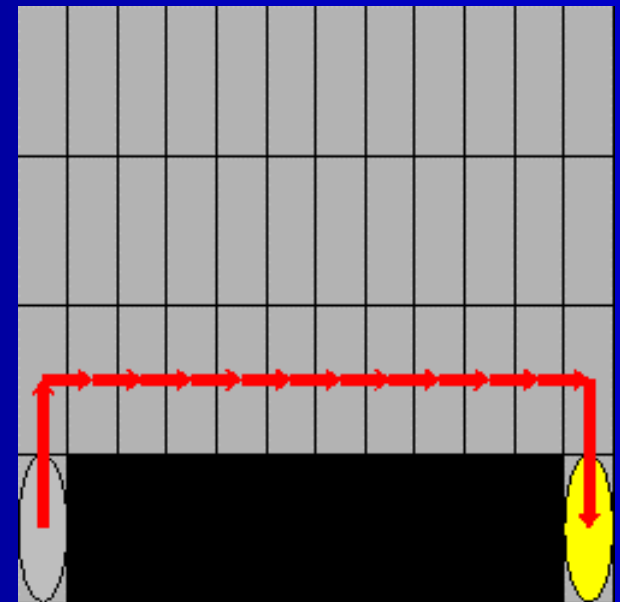
V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference
- VI. Conclusion et Perspectives

- Exemple 2 du *Windy Grid*
– Cliff Walking



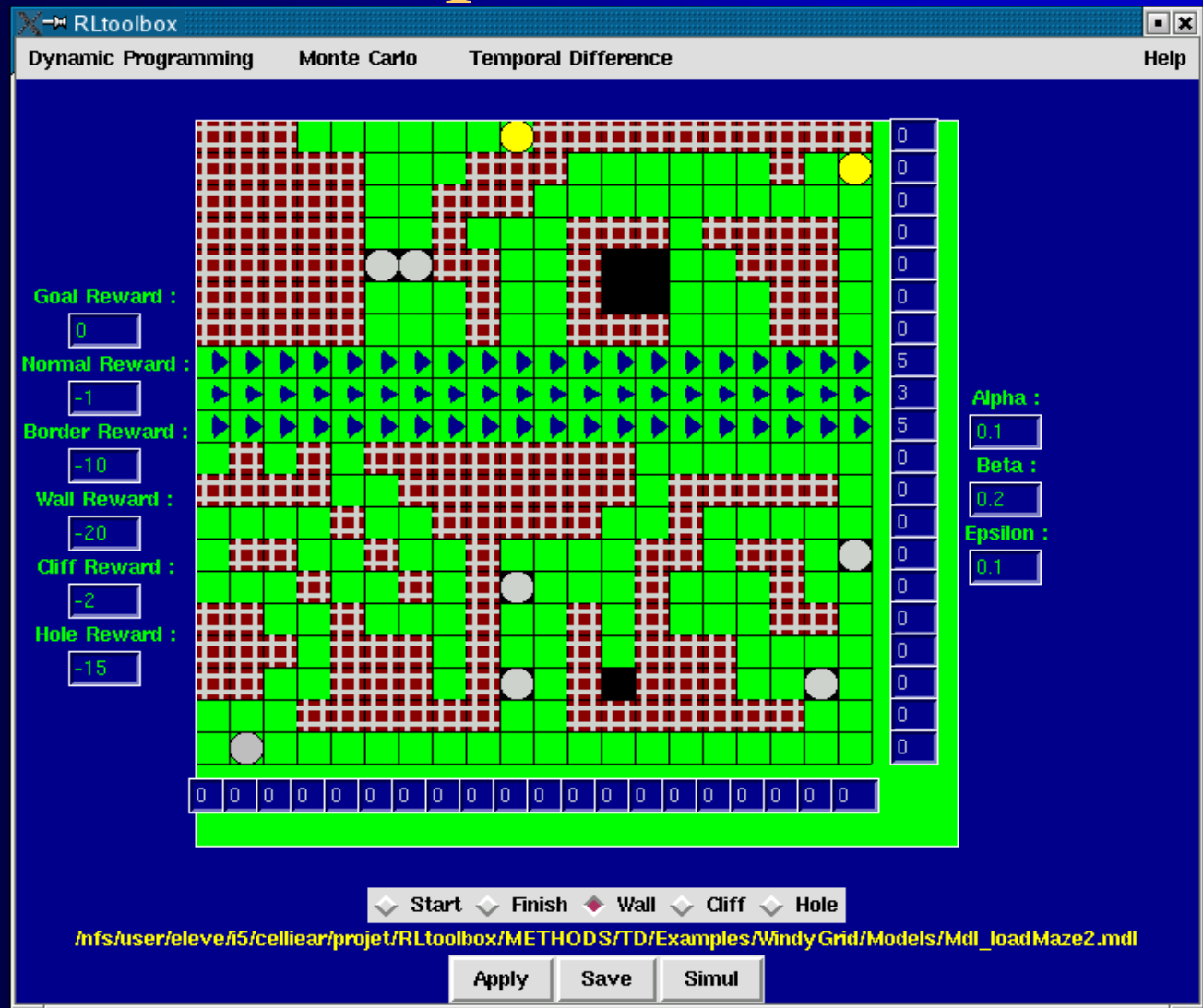
TD Sarsa
500 itérations



TD Q Learning
300 itérations

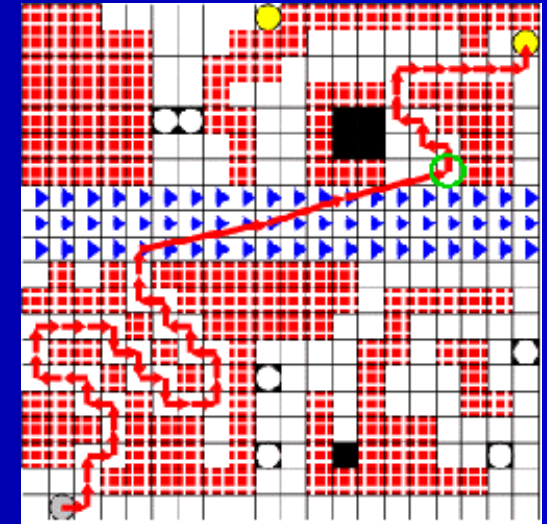
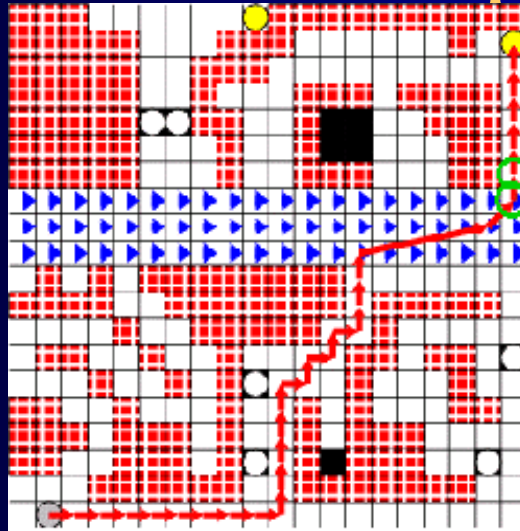
V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference**
- VI. Conclusion et Perspectives

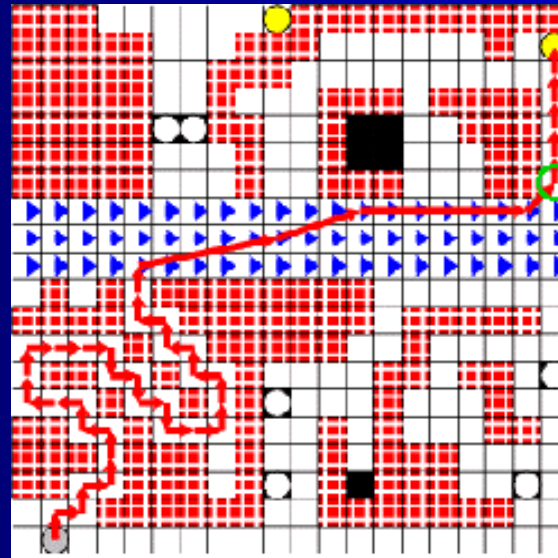


V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. **Temporal Difference**
- VI. Conclusion et Perspectives



Mur : 25
Bord : 25



Mur : -25
Bord : -10

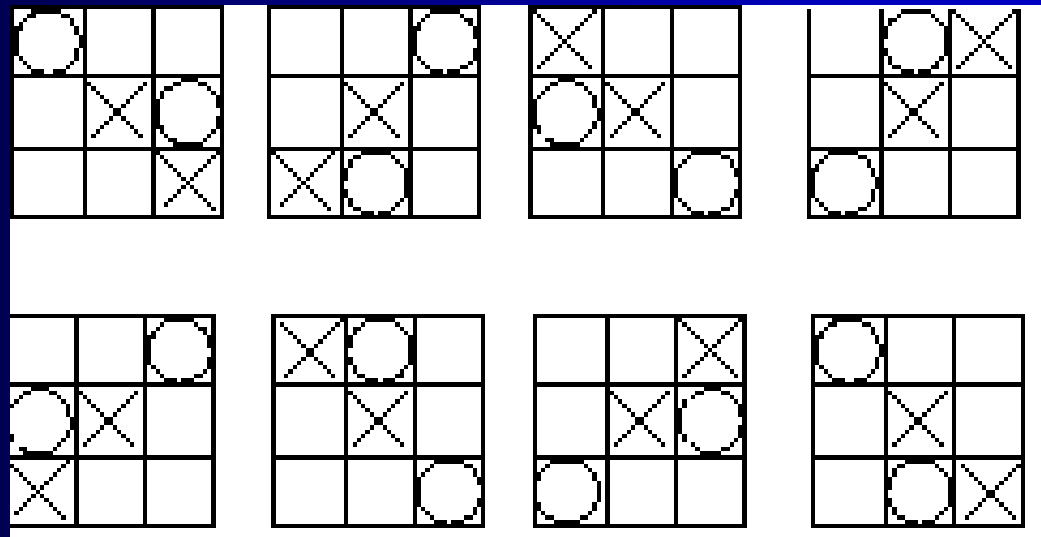
Mur : -10
Bord : -25

V. Temporal Difference

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. Temporal Difference**
- VI. Conclusion et Perspectives

Tic Tac Toe

Compression des états : 5890 \rightarrow 825



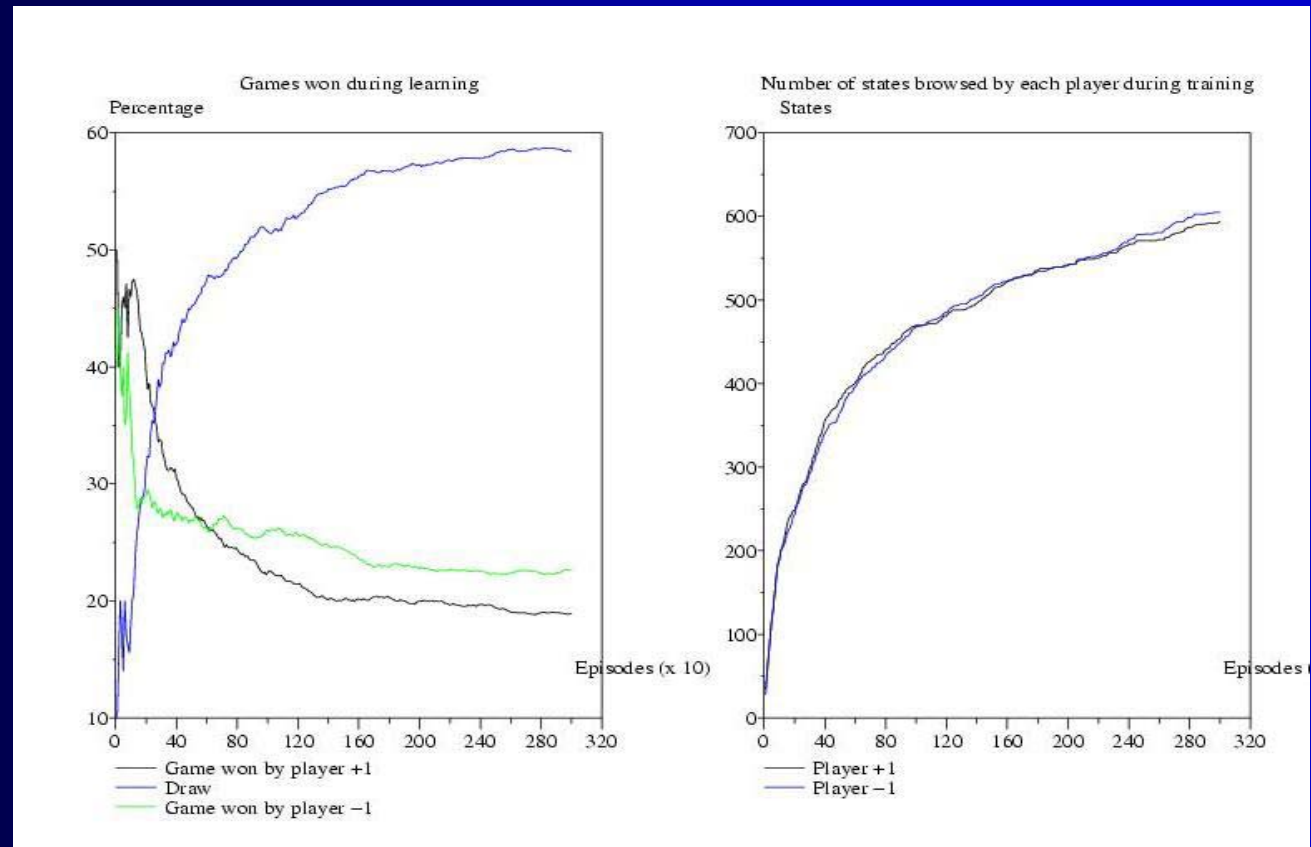
$$S(X) = \sum_{i=1}^9 0.5^i \cdot I_X(i)$$

$$S(O) = \sum_{i=1}^9 0.5^i \cdot I_O(i)$$

V. Temporal Difference

Tic Tac Toe

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. **Temporal Difference**
- VI. Conclusion et Perspectives



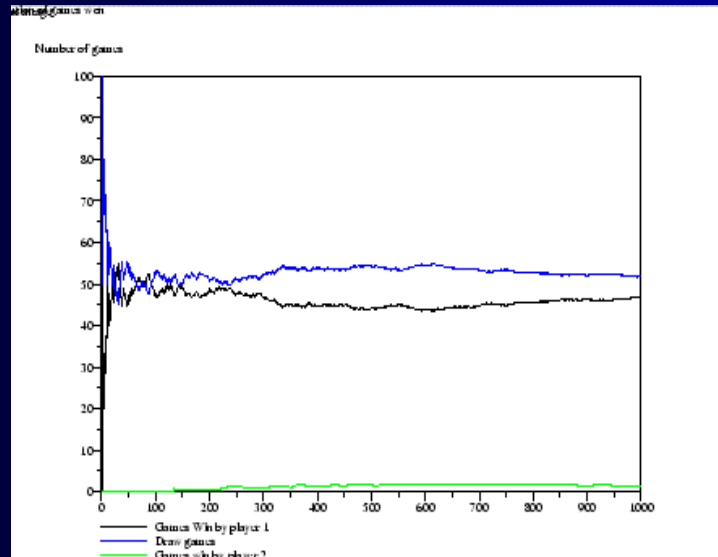
Entraînement sur 3000 parties

V. Temporal Difference

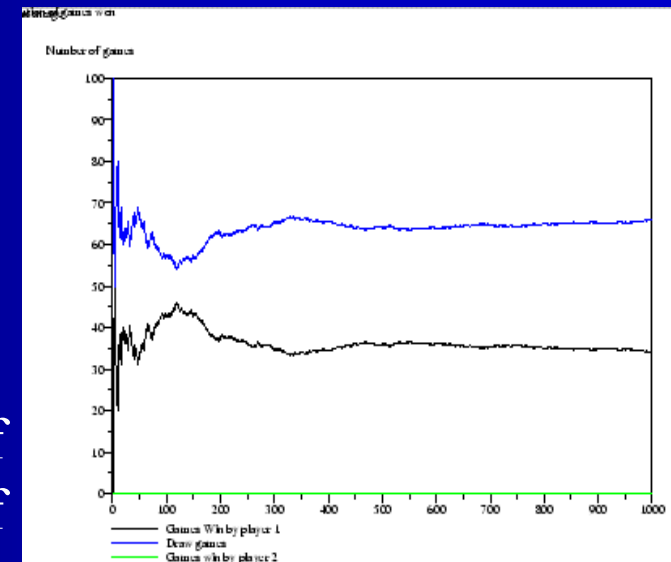
Tic Tac Toe :

Comparaisons d'agents

- I. Introduction
- II. L'apprentissage par renforcement
- III. Programmation dynamique
- IV. Monte Carlo
- V. **Temporal Difference**
- VI. Conclusion et Perspectives



Joueur 1 : 10^4 parties, offensif
Joueur 2 : 10^2 parties, défensif



Joueur 1 : 10^4 parties, défensif
Joueur 2 : 10^2 parties, défensif