

# Programmation Dynamique, Méthodes TD et Optimalité de Bellman

## Introduction à la Programmation Dynamique (PD) et à l'Apprentissage par Renforcement (RL)

La **programmation dynamique (PD)** est une méthode fondamentale pour résoudre des problèmes de prise de décision séquentielle, où les décisions à chaque étape influencent les choix futurs. Ces techniques sont essentielles au développement des **algorithmes d'apprentissage par renforcement (RL)**, utilisés dans des applications comme la planification robotique, les jeux stratégiques et les systèmes de recommandation.

L'objectif de l'apprentissage par renforcement est de permettre à un agent d'apprendre en interagissant avec son environnement, en utilisant des **récompenses** pour guider ses décisions. Les méthodes de PD, initiées par Richard Bellman dans les années 1950, posent les bases théoriques pour une prise de décision optimale dans des environnements dynamiques et incertains.

## Historique et Principes de Base de la Programmation Dynamique

### Origine et Principe d'Optimalité

La **programmation dynamique** a été introduite par **Richard Bellman** avec le **principe d'optimalité** : une politique optimale conserve cette propriété, quels que soient l'état initial et les décisions prises. Ce principe permet de diviser un problème complexe en sous-problèmes plus simples. Bellman a formulé une équation fondamentale, l'**équation d'optimalité de Bellman**, qui exprime la valeur optimale d'un état comme le maximum des récompenses futures attendues.

### Processus de Décision de Markov (PDM)

Les **processus de décision de Markov (PDM)** modélisent la prise de décision dans des environnements incertains. Un PDM est défini par un ensemble d'états, d'actions, de probabilités de transition et de récompenses. Il s'exprime ainsi :

- Ensemble des états  $S$
- Ensemble des actions  $A$
- Fonction de transition  $P(s'|s, a)$  : probabilité de passer de  $s$  à  $s'$  après l'action  $a$

- Fonction de récompense  $R(s, a)$  : récompense pour une action donnée dans un état donné

## Itération de la Politique

L'itération de la politique consiste en deux étapes :

1. **Évaluation de la politique** : Calcul de la fonction de valeur associée à une politique donnée, en itérant jusqu'à la convergence.
2. **Amélioration de la politique** : Ajustement de la politique en choisissant des actions qui maximisent la valeur actuelle de chaque état.

Ce processus converge vers une politique optimale si les étapes sont répétées jusqu'à stabilité.

## Itération de la Valeur

Dans l'**itération de la valeur**, la fonction de valeur est mise à jour itérativement pour approcher la valeur optimale, sans politique intermédiaire. L'algorithme se base directement sur l'équation de Bellman et converge vers la politique optimale une fois la stabilité atteinte.

## Extensions et Applications Avancées de la Programmation Dynamique

### Extensions et Applications de la PD

Durant les années 1970-1980, les méthodes de **DP approximatif** sont apparues pour traiter les problèmes de grande envergure, utilisant des **approximateurs de fonction** (réseaux neuronaux ou fonctions linéaires) pour représenter la fonction de valeur.

### Synergie entre RL et PD

Dans les années 1980-1990, de nombreux algorithmes RL ont été identifiés comme des approximations de méthodes PD. Par exemple, l'**apprentissage par différence temporelle (TD)** combine des concepts de PD et Monte Carlo pour mettre à jour les valeurs d'états en temps réel.

### Développements Modernes en Apprentissage Profond et RL

Avec l'avènement des **réseaux de neurones profonds**, des méthodes de **Deep Reinforcement Learning (DRL)** telles que les réseaux Q profonds (DQN) et les gradients de politique déterministes profonds (DDPG) permettent de traiter des espaces d'états complexes.

# Méthodologie de Résolution des Problèmes d'Apprentissage par Renforcement avec PD

Voici les étapes pour résoudre des problèmes de RL avec PD :

1. **Définition du Problème** : Spécifiez les états, actions, récompenses, et transitions, et modélisez le problème en tant que PDM.
2. **Initialisation** : Initialisez la fonction de valeur  $V(s)$  ou la fonction d'action-valeur  $Q(s, a)$  pour chaque état.
3. **Évaluation de la Politique** : Évaluez la politique avec l'équation de Bellman jusqu'à convergence.
4. **Amélioration de la Politique** : Mettez à jour la politique en choisissant des actions qui maximisent la fonction de valeur.
5. **Extraction de la Politique** (Facultatif) : Une fois la convergence atteinte, extrayez la politique optimale.
6. **Exécution de la Politique** : Implémentez la politique dans l'environnement et évaluez ses performances.

## Applications Pratiques de la Programmation Dynamique en RL

- **Optimisation de la gestion des stocks** : Calcul des quantités optimales à commander pour minimiser les coûts.
- **Planification de la trajectoire en robotique** : Calcul de la trajectoire optimale pour minimiser l'énergie.
- **IA dans les jeux vidéo** : Utilisation des techniques PD pour des décisions optimales des PNJ.

## Conclusion

La programmation dynamique et l'apprentissage par renforcement offrent des solutions puissantes pour les problèmes de prise de décision dans des environnements incertains et dynamiques. Les avancées en **apprentissage profond** permettent aujourd'hui d'appliquer ces méthodes à des environnements avec des espaces d'état de haute dimension, ouvrant des applications dans des domaines tels que la robotique, la finance et la conduite autonome.