

Programmation Dynamique et Apprentissage par Renforcement

Introduction

La **programmation dynamique** est une méthode pour résoudre des problèmes complexes en les décomposant en sous-problèmes plus simples, tout en mémorisant les solutions intermédiaires pour éviter les calculs redondants. Cette approche est utile pour les problèmes de décision séquentielle, où l'objectif est de maximiser une récompense cumulative.

En **apprentissage par renforcement (RL)**, on utilise ces concepts pour optimiser les actions d'un agent dans un environnement, de manière à maximiser les récompenses obtenues à long terme. Ces techniques, y compris le **Q-learning** et le **TD Learning**, reposent sur des équations de récurrence, comme celles introduites par les **équations de Bellman**.

1. Principes de la Programmation Dynamique

La programmation dynamique repose sur deux concepts clés :

- **Sous-problèmes qui se répètent** : Le problème peut être divisé en sous-problèmes similaires. Par exemple, dans la suite de Fibonacci, chaque valeur est calculée à partir des valeurs précédentes.
- **Mémoire pour éviter les recalculs** : Les solutions des sous-problèmes sont mémorisées, ce qui permet de réutiliser ces valeurs sans avoir à les recalculer. Ce processus est souvent appelé ****mémoïsation****.

2. Approches de Programmation Dynamique

Les deux approches pour la programmation dynamique sont :

- **Top-down (mémoïsation)** : Résout le problème principal en le divisant en sous-problèmes de manière récursive, en stockant les résultats pour éviter les recalculs.
- **Bottom-up** : Résout d'abord les sous-problèmes les plus simples, en utilisant les résultats pour résoudre des sous-problèmes de plus grande taille jusqu'à obtenir la solution finale.

3. Exemples de Programmation Dynamique

Exemple : Suite de Fibonacci

Le calcul de la suite de Fibonacci est défini par :

$$F(n) = F(n-1) + F(n-2)$$

En utilisant la programmation dynamique, on stocke les valeurs calculées pour éviter les calculs redondants, ce qui améliore considérablement l'efficacité.

Exemple : Problème du Sac à Dos

Dans le **problème du sac à dos**, on doit maximiser la valeur totale des objets dans un sac à dos de capacité limitée. Si chaque objet i a une valeur v_i et un poids w_i , la solution optimale est donnée par :

$$V(W) = \max(V(W - w_i) + v_i)$$

4. Programmation Dynamique et Apprentissage par Renforcement

En apprentissage par renforcement, la programmation dynamique est utilisée pour optimiser les décisions d'un agent. Les méthodes telles que le Q-learning et le TD Learning reposent sur des **équations récursives** pour estimer les valeurs des actions ou des états en fonction des récompenses futures.

Q-learning et TD Learning

- **Q-learning** : Le Q-learning est une méthode d'apprentissage par différence temporelle (**TD Learning**) qui met à jour une fonction de valeur d'état-action $Q(s, a)$ selon l'équation :

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') \right]$$

où r est la récompense immédiate, γ est le facteur de réduction, et α est le taux d'apprentissage. Cette équation utilise une approximation des récompenses futures pour mettre à jour $Q(s, a)$, optimisant ainsi les décisions sans nécessiter de modèle complet de l'environnement.

- **TD Learning** : Cette famille de méthodes inclut le Q-learning et SARSA. Les méthodes TD combinent la mise à jour incrémentale avec l'apprentissage direct basé sur l'expérience.

Équations de Bellman

Les **équations de Bellman** sont à la base de la programmation dynamique dans les MDP (Processus de Décision Markovien). Elles définissent une relation récursive pour la valeur d'un état ou d'une action :

$$V(s) = \max_a \left[r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \right]$$

et pour la valeur de l'état-action $Q(s, a)$:

$$Q(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$$

Ces équations montrent que chaque valeur dépend des valeurs futures, un principe essentiel de la programmation dynamique.

5. Prédiction de Politique avec Monte Carlo

La méthode Monte Carlo est utilisée pour estimer la fonction de valeur d'une politique donnée en prenant en compte les récompenses futures.

1. Fonction de Retour

Le retour, noté G_t , représente la somme des récompenses pondérées par un facteur de réduction γ :

$$G_t = \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+k+1}$$

2. Estimation avec Monte Carlo de Première Visite

En utilisant Monte Carlo de Première Visite, la valeur d'un état $V(s)$ est calculée en prenant la moyenne des retours pour chaque première visite de cet état :

$$V(s) = \frac{1}{N(s)} \sum_{i=1}^{N(s)} G^{(i)}(s)$$

6. Applications Pratiques

- **Optimisation de la gestion des stocks** : Déterminer les quantités optimales de produits à commander pour minimiser les coûts.
- **Planification en robotique** : Calculer la trajectoire optimale pour minimiser la consommation d'énergie.
- **Jeux vidéo** : Permet aux personnages contrôlés par l'ordinateur de prendre des décisions optimales.

Conclusion

La programmation dynamique et les techniques d'apprentissage par renforcement fournissent des outils puissants pour résoudre des problèmes de décision dans des environnements incertains. En combinant des méthodes comme le Q-learning, TD Learning, et les équations de Bellman, on peut optimiser les décisions d'un agent sur le long terme.

Annexe

Question

Expliquez comment le **Q-learning** utilise les principes de la **programmation dynamique** pour optimiser les actions d'un agent. Dans votre réponse, détaillez l'équation de mise à jour de $Q(s, a)$ en précisant le rôle de chaque terme (récompense immédiate, récompense future, et taux d'apprentissage), et comparez cette approche à celle des **équations de Bellman** dans un MDP (Processus de Décision Markovien).

Réponse

Oui, Q-learning, TD Learning (Temporal Difference Learning), et les équations de Bellman sont des techniques de **Programmation dynamique et apprentissage par renforcement**.

La programmation dynamique est un cadre général pour résoudre des problèmes de décision séquentiel.

En apprentissage par renforcement, on utilise ces techniques pour optimiser les actions d'un agent dans un environnement.

Q-learning et TD Learning

1. Q-learning :

C'est une méthode de Temporal Difference (TD) Learning. Elle met à jour une fonction de valeur d'état-action.

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') \right]$$

2. Equations de Bellman :

Les équations de Bellman définissent une relation de récursivité pour la valeur d'un état ou d'une action.

$$V(s) = \max_a \left[r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \right]$$
$$Q(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$$

Pourquoi ces méthodes relèvent de la programmation dynamique

Elles utilisent des équations récursives pour estimer les valeurs en fonction des valeurs futures, en mémorisant les résultats.