

NodeP - free RL

NodeP-free → sans modèle → Des situations où l'agent n'a pas des connaissances de proba de transition entre les états et des récompenses associées. + ces d'actions discrètes → Modélisation par MDP (Markov)

Commencer sans info → explorer l'env pour acquies ces informations

Agent: teste chaque état et chaque transition pour découvrir les récompenses.

se fait plusieurs fois pour pouvoir estimer les proba de transitions

Q-Learning: NodeP-free (sans modèle)

apprend directement la fonction Q à partir de l'expérience sans avoir besoin d'un modèle de l'env

À l'initialement: l'agent connaît les états et les act° possibles

Approche du Non-Case pour estimer Q:

→ l'agent apprend à partir d'épisodes complets d'interactions avec l'env.

→ ne attend la fin de l'épisode pour mettre à jour les valeurs des états.

→ Après chaque épisode, l'agent attribue des récompenses réelles à chaque état visité pour obtenir une estimation de la valeur de chaque état.

→ mise à jour de la valeur: en fonction des récompenses réelles obtenues.

Episode: état visité → mise à jour?

$$V_{\text{new}} = \text{val} + \frac{1}{n} (R - \text{val})$$

TD Learning:

→ Apprentissage par différence temporelle

→ met à jour les valeurs des états à chaque pas de temps: l'agent apprend plus rapidement et de manière incrémentielle.
→ Estimation de la valeur de l'état actuel: en utilisant R, la val de l'état suiv pondue par un facteur d'apprentissage.

Exploration / Exploitation

Code: Greedy

class EpsilonGreedyPolicy:

def __init__(self, epsilon):

self.epsilon = epsilon

def select_action(self, q_values):

p = random.uniform(0, 1)

if p < self.epsilon:

action = random.choice(range(len(q_values)))

else:

action = max(range(len(q_values)), key =

lambda a: q_values[a])

return action

Epsilon_value = 0.1

Epsilon_greedy_policy = EpsilonGreedyPolicy(Epsilon_value)

q_values_ex = [0.2, 0.8, 0.1, 0.3]

select_act = Epsilon_greedy_policy.select_action(q_values_ex)

Choix Epsilon → une instance de classe → sélection d'une act° en fonction de Q-values

remplace par les Q-val de l'env

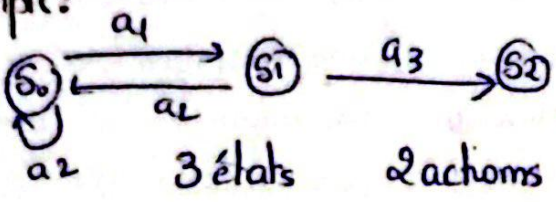
SARSA:

→ C'est un algorithme qui met à jour les valeurs Q en utilisant une stratégie de politique ϵ -greedy. Voir le code

Q-Learning: env + Act° → Q-learning

→ C'est un algo qui apprend les Q-values pour chaque état-action dans un env.

Exemple:



3 états 2 actions
 $R(S_0) = -0.4$ $R(S_1) = -0.1$ $R(S_2) = 1$
 $\alpha = 0.5$ $\gamma = 0.5$

au dem: l'agent en S_0 et choisit a_2

La formule de la mise à jour de la Q-table

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

• (S_0, a_1)
 Act² à prendre $\pi(S_0) = \argmax(Q(S_0, a_1), Q(S_0, a_2))$
 $= \argmax(0, 0) = a_2$

• $(S_0, a_2) \rightarrow (S_0, a_1)$
 $Q(S_0, a_2) = Q(S_0, a_2) + \alpha (R(S_0, a_2) + \gamma \max_{a'} \{Q(S_0, a_1), Q(S_0, a_2)\} - Q(S_0, a_2))$
 $= 0 + 0.5 (-0.1 + 0.5 \max\{0, 0\} - 0) = -0.05$

$\pi(S_0) = \argmax \{Q(S_0, a_1), Q(S_0, a_2)\} = a_1$

• $(S_0, a_1) \rightarrow (S_0, a_2) \rightarrow (S_1, a_1)$
 $Q(S_0, a_1) = -0.05$
 $\pi(S_0) = a_2$

• $(S_0, a_1) \rightarrow (S_0, a_2) \rightarrow (S_0, a_1) \rightarrow (S_1, a_2) \rightarrow (S_0, a_1)$
 $Q(S_1, a_2) = -0.0625$
 $\pi(S_0) = a_1$

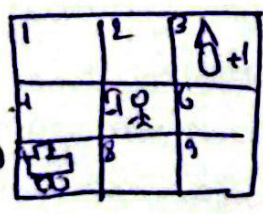
• $(S_0, a_1) \rightarrow (S_0, a_2) \rightarrow (S_0, a_1) \rightarrow (S_1, a_2) \rightarrow (S_0, a_1) \rightarrow (S_1, a_1)$
 $Q(S_0, a_1) = -0.075$
 $\pi(S_1) = a_3 \rightarrow (S_2)$ Etat terminal
 $Q(S_2, \text{Nul}) = 1$
 $Q(S_1, a_3) = 0.2$

Etat	Act ²	
S_0	$Q(S_0, a_1) = -0.075$	$Q(S_0, a_2) = -0.05$
S_1	$Q(S_1, a_2) = -0.0625$	$Q(S_1, a_3) = 0.2$

Exemples:

• l'agent initialement à 7

Les déplacements: droit, gauche, haut, bas



Les récompenses: -1 en 5, 1 en 3, 0 partout

→ Epoches = 100 $\alpha = 0.1$ $\gamma = 0.9$

Etat 1 → 9 $Q = 10$ 9x4

0 → 1 (↑) 0 → 1 (↓) 0 → 2 (←) 0 → 3 (→)

La fonction take-action détermine l'action adoptée:

```

def take-action(st, Q, eps):
    if random.uniform(0,1) < eps:
        action = randint(0,3)
    else:
        action = mp.argmax(Q[st])
  
```

return action

Mise à jour de la Q-table:

```

while not env.is_finished():
    at = take-action(st, Q, 0.4)
    stp1, r = env.step(at)
    alp1 = take-action(stp1, Q, 0.0)
    Q[st][at] = Q[st][at] + 0.1 * (r + 0.9 *
    Q[stp1][alp1] - Q[st][at])
    st = stp1
  
```

• Deep Q-Learning:

→ utiliser les réseaux de neurones profonds pour représenter et estimer la fonct² Q

→ Architecture:

• Estimer la cible y_i et $Q(s,a)$ à l'aide de deux réseaux neuronaux distincts.

Target Network Q-Network Separat²
 • θ_{i-1} : (poids, biais) ∈ Réseau cible des réseaux
 θ_i ∈ Réseau Q
 • s: état actuel

réseau Q: calcul les val d'act² $Q(s,a)$

réseau cible: calcul les $Q(s',a)$ en utilisant l'état suivant pour comb. la cible TD.

Experience Replay:

→ Utiliser pour améliorer l'efficacité et la stabilité de l'entraînement: l'agent stocke les expériences dans une mémoire de répétition et échantillonne un lot d'exp de manière aléatoire pour caractériser l'entraînement.

Day 1 recap :

→ Structure de données utilisée pour stocker les expériences passées en vue de l'entrainement.

→ Tableau qui contient un nombre fixe des expériences les plus récentes

→ Chaque exp est représentée par un tuple

(s, a, r, s') → l'état suivant
état actuel → action effectuée → récompense obtenue

Deep reinforcement Learning

Actor Critic

Le critique : méthode qui évalue chaque action que l'agent prend dans l'environnement avec une scalaire > 0 ou < 0 (à l'échelle DDN).

Acteur : politique paramétrée qui définit comment les actions sont sélectionnées

Policy Gradient : paramètre la fonction de politique et met à jour les paramètres dans la direction qui augmente le rendement attendu.

→ Le critique évalue les actions choisies par l'acteur

→ L'acteur est mis à jour dans la direction suggérée par le critique ce qui stabilise l'entraînement et réduit la variance

Policy Gradient Methods : optimisent directement la politique

Actor-Critic Methods : utilisent une fonction de valeur pour guider la mise à jour de la politique.

Politique : stratégie qui décide comment un agent devrait prendre des décisions dans un environnement donné.

DDO algorithms

→ appartient à la famille Policy Gradient

→ Résolve certains problèmes de stabilité et de convergence associés à d'autres algo de Policy Gradient.

- initialise un réseau neuronal pour représenter la politique, $\pi_\theta(a/s)$ où θ sont les paramètres du réseau

- initialise un réseau neuronal pour représenter la valeur $V_\phi(s)$ où ϕ sont les paramètres du réseau

1 énoncé

- interagit avec l'environnement pour collecter des états / actions / récompenses
- Calculez les avantages $A(s, a)$ en utilisant les récompenses observées et les estimations
- Calculez la fonction objective $P_{\theta} = \text{combiné de l'obj politique } L^{\text{clip}}$ et l'obj de valeur L^{vf}
 - ↳ maximiser la proba d'actions qui ont été prises
 - ↳ minimiser la diff quadratique moyenne entre les estim et les avantages calculés
- utiliser un algo d'optimisation pour mettre à jour les params du réseau de politique et du réseau de valeur en minimisant P_{θ} .

$$\pi(\theta) = \frac{\pi_{\theta}(a_t/s_t)}{\pi_{\theta_{\text{old}}}(a_t/s_t)}$$

> 1
more likely on the current policy than the old policy

$0 < \cdot < 1$
less likely for the current policy than for the old one

The objective functions:

$$L^{\text{clip}}_{\theta} = \mathbb{E} \left(\min \left(\pi(\theta) \cdot \hat{A}_t, \text{clip}(\pi(\theta), 1-\epsilon, 1+\epsilon) \cdot \hat{A}_t \right) \right)$$