

Introduction au renforcement Learning

Machine Learning: apprentissage non supervisé:

Obj: découvrir les structures sous-jacentes à des données non étiquetées

Données: Données non étiquetées

Exp: détection d'anomalie

Machine Learning: Apprentissage supervisé

Obj: apprendre des fonctions de prédiction à partir des données annotées

Données: Données étiquetées

Exp: prédiction d'une valeur dans des problèmes de régression classification

Reinforcement Learning: Apprentissage par renforcement

Obj: Laisser l'algo apprendre de ses propres erreurs.

choix aléa au début \rightarrow se trompe \rightarrow puni

\hookrightarrow sinon \rightarrow récompense

Données: Etats et Actions

Obj de l'app supervisé et non supervisé:

\hookrightarrow La minimisation de l'erreur

Obj de l'apprentissage par renforcement:

\hookrightarrow maximiser la récompense

Les principaux composants de l'app par renforcement:

• Agent: l'entité qui prend des décisions dans l'environnement.

• Etat: Représentation de la situation actuelle de l'agent dans l'env. $S = \{ \text{position, vitesse, ...} \}$

• Action: La décision prise par l'agent à un moment donné $A = \{ \text{Act1, ...} \}$

• Récompense: $\forall t, \forall st, \forall at$, il reçoit une

bien et si \hookrightarrow sinon.

En RL: un agent procède à des observations et réalise des actions au sein d'un env. Il reçoit en retour des récompenses.

Exemple:

Finance:

Agent: observe les prix des actions

Décide le nbr d'act à acheter/vendre

Récompense: dépendent des gains et des pertes

Politique: stratégie à suivre pour atteindre un obj

$\pi: S \rightarrow A$: associe à chaque état s une action a

Exemple: Aspirateur

\rightarrow Récomp: volume de pouss ramassé en 30min

\rightarrow Politique: avancer avec proba $= p$
tourner avec proba $= 1-p$

Exemple: Cart Pole

\rightarrow Politique: déclenche une accélération vers la gauche lorsque le bâton penche vers la gauche et vice versa.

Evaluation des actions:

Dans $st \in S$, l'agent envoie l'action $at \in A$ et reçoit une récompense r_{t+1}

seq de récompense: r_{t+1}, r_{t+2}, \dots

\rightarrow la fonction de récompense:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

$\gamma \in [0, 1]$: paramètre de rabais

\hookrightarrow détermine la valeur présente des récompenses futures

$\gamma = 0 \Rightarrow$ agent myope: maximiser les récompenses immédiates

$\gamma \rightarrow 1 \Rightarrow$ agent avec un horizon de plus en plus lointain

OpenAI Gym 1: c'est un défis qui a un environnement fonctionnel pour

entraîner un agent, il faut lui donner env simulés tq les jeux Atari, des jeux de société, des simulations physiques 2D ou 3D qui permet de entraîner, comparer et développer des algorithmes d'apprentissage par renforcement

Cartpole environment I

On utilise pour tester et développer des algo
L'obj: maintenir un poteau en équilibre aussi longtemps que possible en déplaçant un chariot horizontalement.

import gym
env = gym.make('CartPole-v1')
obs = env.reset()
obs → array([-0.0127, -0.00156, 0.042, -0.0018])



gym.make: créer un environnement
env.reset(): réinitialiser l'env, et renvoie une obs initiale
action → réinitialiser

→ Les obs dépendent du type d'env

L'env: Cartpole, alors les obs sont les réels:
La position horizontale du chariot, (entre 0,0)
sa vitesse (≥ 0 vers la droite, ≤ 0 vers la gauche)

L'angle du poteau (0,0 vertical)

Vitesse Ang du poteau ⊕ Angle orienté

env.render(): affiche une fenêtre contextuelle à chaque nouvelle action prise à chaque étape.

env.render() → affiche l'environnement
True

env.action-space: pour définir les caractéristiques de l'espace d'action de l'env: on peut définir si l'espace d'act° est continu ou discret, définir les valeurs min et max des actions, ...

env.action-space

→ Discrete (2) → c'est les actions possibles sont 0 et 1: l'accélération vers la gauche (0) ou vers la droite (1).

→ peut avoir des autres act° supp pour les autres env, ou d'autres type d'action

env.observation-space

→ Box (4,) vecteur à 4 dimensions réelles

env.observation-space.low } découvrant

env.observation-space.high } la plage

de chaque variable d'observation

Les valeurs min et max

env.step(): effectuer une act° à chaque étape. L'act° est spécifiée en tant que paramètre

env.step(): retourner le paramètre (obs, récompense, done et info)

Le poteau → chute alors obs (2) = vitesse

L'acc est vers la droite, action = 1

obs, reward, done, info = env.step(action)

[...] recomp boolean → indique si l'env a fini
si l'env a fini de réinitialiser l'env à l'état initial
de l'info diagnostique

Ecrire un Script qui teste une politique afeet dans l'env Cartpole

import gym

import random

env = gym.make('CartPole-v1')

def test_policy(env, nb_ep=5):

for e in range(nb_ep):

obs = env.reset()

total_reward = 0

while True:

env.render()

act = env.action_space.sample()

obs, reward, done, info = env.step(action)

total_reward += reward

if done:
print('épis:', e+1, 'Rec:', total rew)
env.close()

Processus de Décision de Markov:

• hypothèse de Markov: $\forall t$, la proba de passer de s à s' en faisant une act° a
 $P(s_{t+1} = s' | s_t = s, a_t = a) \rightarrow$
 ne dépend pas des états précédents

Un problème de RL est définie comme un PDM (chaînu de Markov) où T et R sont inconnus

$T: S \times A \times S \rightarrow [0,1]$: fonct° de transition qui donne la proba d'arriver en s'

$$T(s, a, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$$

$R: S \times A \rightarrow \mathbb{R}$: fonct° de récompense

$$R(s, a) = E(R_t | s_t = s, a_t = a)$$

Fonction valeur d'un état s :

$$V^\pi(s) = E(R_t | s_t = s)$$

$$= E(\sum \gamma^k R_{t+k+1} | s_t = s)$$

Fonction état-valeur:

$$Q^\pi(s, a) = E(R_t | s_t = s, a_t = a)$$

$$= E(\sum \gamma^k R_{t+k+1} | s_t = s, a_t = a)$$

Fonction valeur optimale

$$V^*(s) = \max_a [R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s')]$$

→ Si l'agent agit de façon optimale

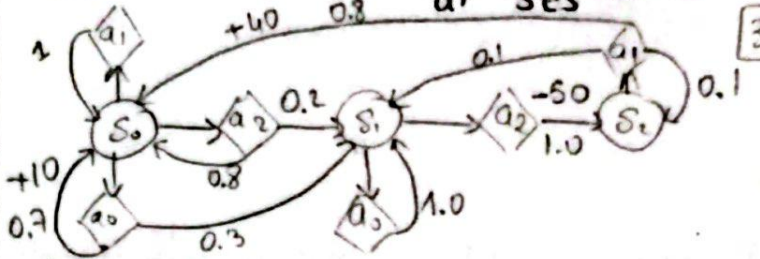
$$V_{opt}(\text{état}_{\text{current}}) = R(a_{optim}) + \sum V_{opt}(\text{état}_{\text{suiv}})$$

$$V_{k+1}(s) = \max_a [R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V_k(s')]$$

Equation d'optimalité de Bellman pour estimer la fonct° état-valeur optimale

$$Q^*(s, a) = R(s, a) + \gamma \max_{a'} [\sum_{s' \in S} T(s, a, s') Q^*(s', a')]$$

$$Q_{k+1}(s, a) = R(s, a) + \gamma \max_{a'} [\sum_{s' \in S} T(s, a, s') Q_k(s', a')]$$



3 états, 3 actions,

$S_0 \rightarrow$ choisir $a_1 \rightarrow$ reste dans S_0 avec certitude et sans aucune récompense
 \rightarrow choisir $a_0 \rightarrow$ 70% chance d'obtenir une récompense et reste dans S_0
 \rightarrow choisir $a_2 \rightarrow$ retour dans S_1 avec proba de 0.2
 reste immobile $\leftarrow a_0 / a_2 \rightarrow -50$ jusqu'à $a_1 \leftarrow S_2$

Q-value iteration algo

La c'est une méthode utilisée pour déterminer la fonction d'action-valeur optimale d'un proc de décision markovien

Q : la récompense cumulée attendue de prendre une action spécifique dans un état particulier.

1/ Définir le PDP

→ transitions-proba
 → rewards
 → possible actions

2/ $Q_{k+1}(s, a) \quad \forall (s, a)$

3/ gamma = 0.9

for itération un range(50):

$Q_{\text{prev}} = Q_{\text{value}} \text{ copy}()$

for s in range(3):

for a in possible_actions[s]:

$$Q_{\text{value}}[s, a] = \text{mp.sum}([\text{transition-proba}[s][a][sp] * (\text{rewards}[s][a][sp] + \text{gamma} * \text{mp.max}(Q_{\text{prev}}[sp]))] \text{ for } sp \text{ in range(3)})]$$

→ Q-values result: array([[1., 1., 1.], [0.17, ..., -inf, ...], [1., 1.]])

Q values: 18,92 17,03 13,62
 0 -∞ -4,88
 -∞ 50,13 -∞

$(s_0) \rightarrow (a_1)$
 +17 ← rewards

Pour chaque état, l'action qui a la Q-valeur la plus élevée:

mp. argmax (Q-values, axis = 1)
 array ([0, 0, 1])

Catégorisation des techniques de RL:

• RL basé sur la valeur:

→ valeur Q: Estimation des Q-valeurs représentant la qualité des actions dans un état donné.

Q-Learning / SARSA

→ valeur V: Estimer la valeur d'un état particulier

Dynamic programming

• RL basé sur la politique:

→ Politique déterministe: qui mappe les états aux actions directement

Policy Gradient Methods

→ Politique stochastique: les politiques sont traitées comme des distributions de proba sur les actions

Trust Region Policy Optimization

• Méthodes basées sur la modélisation:

→ Model-based: utilisent des modèles internes pour simuler l'env et planifier des actions

NPC

→ Model-free: ne nécessitent pas de modèle explicite de l'env.

Q-learning / Policy Gradient Methods

• Exploration et exploitation:

→ Exploration: ^{essayer de nouvelles actions} pour découvrir de nouvelles actions et états pour améliorer la

Connaissance de l'agent sur l'env.

Epsilon Greedy

→ Exploitation: utilisation de la connaissance actuelle pour prendre des décisions optimales

Greedy Policy

bonne
Récompense

• Deep Learning:

→ Réseaux de neurones: intégration de réseaux de neurone pour représenter et apprendre des politiques ou des valeurs plus complexes.

Deep Q-Networks (DQN)

Actor Critic, DDQN

Deep Policy Network, NPC

⇒ Ces techniques sont utilisées pour entraîner des agents intelligents à prendre des décisions dans un environnement pour maximiser une récompense cumulative au fil du temps.