

**Master 2 Analyse et Politique Économique – DS2E
Semestre 1**

UE Machine Learning – Reinforcement Learning

Solving Cliff Walking with Reinforcement Learning

Sous la responsabilité de M. Bertrand Koebel, Professeur de
Sciences Économiques à la Faculté des Sciences Économiques
et de Gestion
Université de Strasbourg

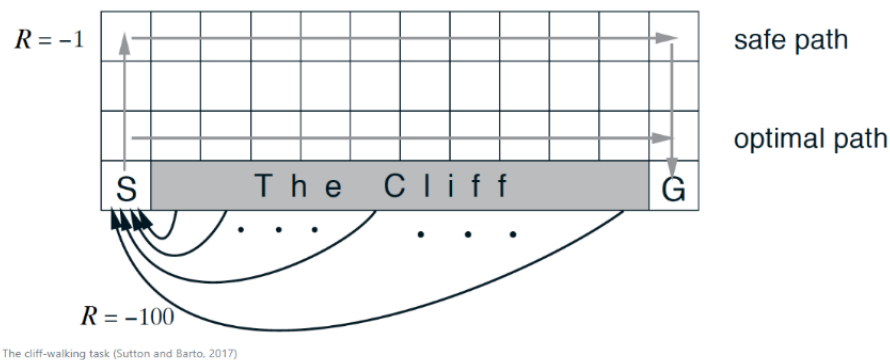
Zachary JANATI, Chafi KHERMOUCHE

1. Description de l'environnement

La "Marche sur les falaises" représente un scénario classique défini par Sutton et Barto en 2018, dans lequel un agent aspire à naviguer depuis la tuile en bas à gauche jusqu'à la tuile en bas à droite. L'objectif de l'agent est de minimiser le nombre de pas tout en évitant la zone périlleuse, symbolisée par la falaise. Un épisode de ce scénario se conclut soit lorsque l'agent marche hors des limites sur la falaise, engendrant une grande récompense négative, soit lorsqu'il atteint la tuile cible, générant une récompense positive.

Les tuiles dans cet environnement peuvent être catégorisées comme suit :

- Tuile sûre : Dans cette configuration, l'agent peut avancer en toute sécurité sans risque de conséquences néfastes.
- Tuile dangereuse (la falaise) : Cette zone constitue un danger potentiel, entraînant une situation où l'agent serait bloqué de manière permanente, mettant ainsi fin au jeu.



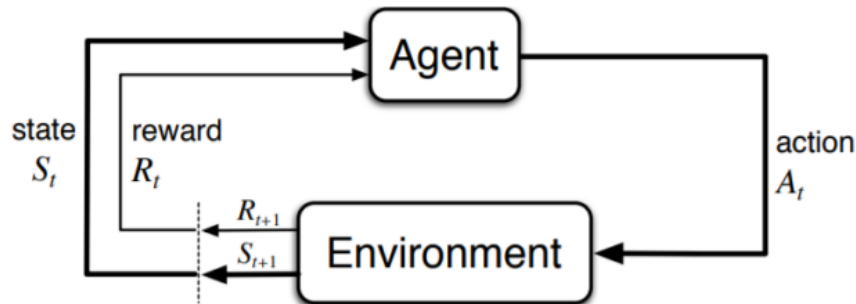
L'environnement est défini par une grille de 4x12 carreaux, où chaque carreau est identifié par une lettre spécifique

L'agent a 4 actions possibles :

- ➔ Aller à GAUCHE
- ➔ Aller en BAS
- ➔ Aller à DROITE
- ➔ Aller en HAUT

2. Interaction entre l'environnement et l'agent (actions)

L'apprentissage par renforcement est un concept central en machine learning, où l'apprentissage se produit à travers un processus itératif de tentatives et d'erreurs. Une manière de visualiser ce concept est à travers un diagramme qui présente deux acteurs principaux : l'agent (comme un marcheur) et l'environnement (comparable à une montagne avec des falaises).



L'agent observe son environnement et utilise ces observations, ainsi que des récompenses le cas échéant, pour prendre des décisions. L'objectif premier de l'agent est de maximiser la récompense totale attendue. Dans ce contexte, imaginons la marche comme un jeu, où l'objectif de l'agent est de passer de l'état de départ (D) à l'état de victoire (V) sans tomber dans les états de piège (T). Chaque état représente une position possible sur la surface, et pour chacun de ces états, l'agent peut choisir parmi quatre actions : aller à gauche, en bas, à droite ou en haut.

L'apprentissage dans ce contexte implique de déterminer quelle action choisir dans chaque état. Pour ce faire, il est utile d'attribuer une valeur de qualité (Q-value) à chaque action pour chaque état. Au total, il y a 48 états (16×4) et 4 actions possibles, ce qui signifie que l'agent doit calculer $48 \times 4 = 192$ valeurs de Q. Pour organiser ces informations, une table peut être utilisée, où les lignes correspondent aux états et les colonnes aux actions. Chaque cellule de cette table contient la valeur $Q(s, a)$, qui représente la qualité de l'action a dans l'état s . Une valeur de Q élevée indique que l'action est plus prometteuse, tandis qu'une valeur de Q basse indique qu'elle est moins recommandée.

Lorsque l'agent se trouve dans un état et doit prendre une décision, il peut consulter cette table pour déterminer quelle action a la valeur de Q la plus élevée. Cependant, il est important de noter que choisir systématiquement l'action avec la valeur de Q la plus élevée n'est pas nécessairement la stratégie optimale. Parfois, l'exploration d'autres actions, même celles avec des valeurs de Q plus basses, peut être bénéfique, car cela peut conduire à de meilleures découvertes à long terme. L'apprentissage par renforcement vise à équilibrer l'exploitation (choisir la meilleure action connue) et l'exploration (essayer de nouvelles actions) pour atteindre l'objectif de maximisation des récompenses cumulées attendues.



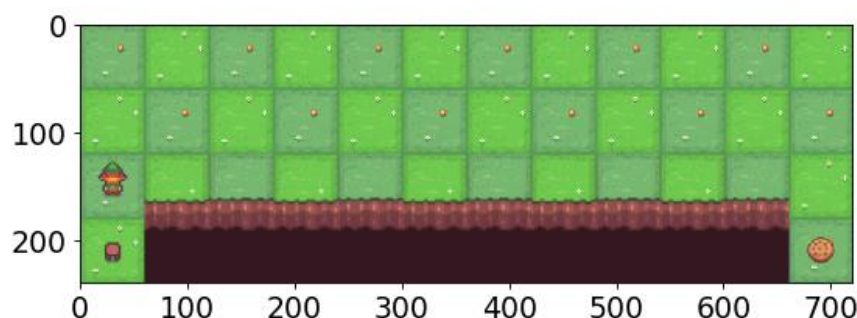
Le tableau ci-dessous est complémentaire et permet de bien visualiser le chemin que l'agent cherchera à prendre, avec deux « chemins » : un chemin risqué et un chemin optimal.

Le chemin optimal est celui qui maximise les récompenses de l'agent tout en minimisant les pénalités, le conduisant de manière sûre et efficace de la position de départ à la position de but sans tomber de la falaise. En revanche, un chemin risqué pourrait comporter des déplacements qui passent près de la falaise ou qui nécessitent des actions plus hasardeuses, augmentant le risque de tomber et de recevoir des pénalités importantes, ce qui se traduit par une série de valeurs d'état plus négatives et une stratégie globalement moins efficace.

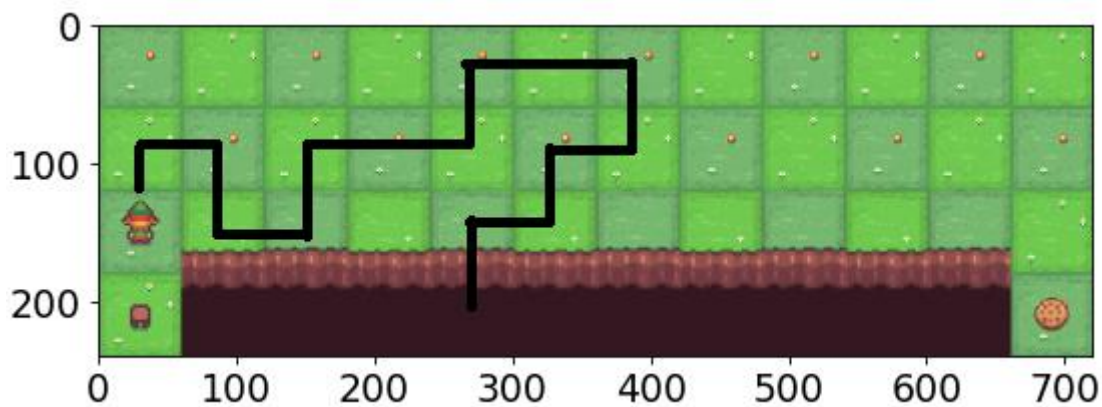
3. Le modèle non déterministe

Lorsque nous envisageons de traduire les actions de chaque état de la table ci-dessus en termes de probabilités, nous constatons que, pour chaque action, il y a une probabilité équivalente de $1/3$ (environ 0.33) de se retrouver dans le carreau souhaité. Cela s'explique par la présence d'une possibilité de glisser et de se retrouver dans un endroit non souhaité, représentant une probabilité de $2/3$ (environ 0.67). Dans de tels cas, on qualifie le modèle de non déterministe, indiquant une incertitude quant à l'issue de l'action entreprise, car il existe une chance non négligeable de se retrouver dans une situation non prévue (probabilité $p=2/3$). Ainsi, Par exemple, l'agent commence sur l'état « Départ » au début, et si l'agent choisit l'action « en haut » puis « à droite », nous aurons les probabilités suivantes dans le modèle non déterministe :

- ➔ $P(\text{En bas}) = 1/4 = P(\text{Perte de la partie})$
- ➔ $1-P(\text{En bas}) = P(\text{Se déplacer ailleurs}) = 3/4$



De manière logique, il est observé qu'une des solutions pour atteindre l'état de "victoire" (V) consiste à effectuer une montée, suivis de dix déplacements vers la droite, et enfin d'une descente (politique 1). Cependant, en raison de la nature non déterministe du modèle, cela entraînera une probabilité de succès, notée $p(V=1|politique\ 1)$, équivalente à $(1/3)^2 \times (1/4)^{10}$ (0.4% ce qui est très faible) , si l'on suit les actions spécifiées précédemment par la politique 1.



Exemple de situation possible dans le cadre du modèle non déterministe

Comme nous pouvons le remarquer ci-dessus dans le modèle non déterministe, l'agent tombe dans la falaise malgré qu'on lui ait indiqué de suivre la procédure par la politique 1 puisque pour une action imposée à l'agent, il a une probabilité plus grande ($3/4$) de se retrouver dans un carreau non souhaité.

4. Le modèle déterministe

Dans un modèle déterministe, les actions prescrites par une politique sont exécutées avec une probabilité de 1, ce qui signifie qu'elles sont certaines. Cette certitude découle du fait que chaque action entreprise par l'agent mène de manière invariable à un état spécifique, conformément aux règles définies par le modèle. Dans le contexte de la politique 1, l'agent suit les directives avec une certitude totale, se déplaçant de manière précise vers les emplacements spécifiés jusqu'à ce qu'il atteigne l'état de victoire ($s=V$).

En d'autres termes, dans un modèle déterministe, l'agent a une connaissance parfaite des conséquences de ses actions, et les résultats sont prédéterminés à partir de l'état initial. Cette prévisibilité renforce la certitude associée à l'exécution des actions prescrites par la politique en question, ce qui contraste avec les modèles non déterministes où l'issue des actions peut varier en raison de l'incertitude introduite par des éléments aléatoires.

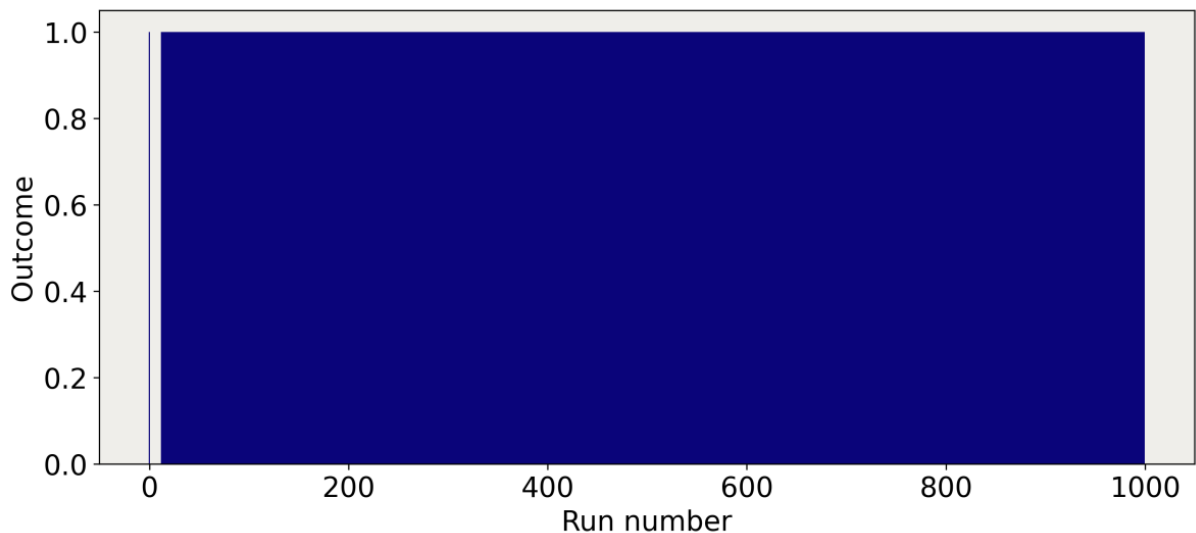
■ Q-LEARNING

Supposons que notre agent a réussi à atteindre l'état de victoire ($s=V$). La question se pose de savoir comment répercuter cette information jusqu'à l'état initial pour améliorer le processus d'apprentissage. C'est précisément ce que le Q-learning propose comme solution à ce défi.

Ainsi, en utilisant le Q-learning, nous mettons à jour les informations stockées dans la table Q pour refléter les meilleures décisions possibles dans chaque état. Cela permet à l'agent d'apprendre progressivement à optimiser ses actions pour maximiser les récompenses à long terme en se basant sur les récompenses instantanées obtenues et les estimations de la valeur future. C'est un mécanisme clé de l'apprentissage par renforcement qui permet à l'agent de s'améliorer au fil du temps en utilisant les informations acquises au cours de ses expériences passées.

Q-table before training:

[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]	[0. 0. 0. 0.]
[0. 0. 0. 0.]			



Nous pouvons observer que l'agent a accompli un apprentissage fructueux très rapidement. Chaque barre bleue dans la représentation graphique symbolise une réussite, illustrant ainsi que l'agent a rencontré des difficultés initiales pour atteindre la victoire au tout début du processus d'entraînement. Toutefois, après avoir expérimenté la réussite de manière répétée, à partir de certaines itérations spécifiques (notamment à partir de ... répétitions), l'agent a démontré une amélioration significative, manifestant une capacité à remporter régulièrement et à chaque fois des victoires.

L'analyse de la table-Q entraînée s'avère tout aussi fascinante. Ces valeurs reflètent la séquence unique d'actions que l'agent a assimilée pour atteindre son objectif. Cette représentation offre un aperçu détaillé des connaissances acquises par l'agent au cours de son processus d'apprentissage, décrivant les choix d'actions spécifiques qui ont conduit aux succès observés dans la tâche.

b) Modèle non déterministe

Dans le cas du modèle non déterministe, nous pouvons voir qu'il est quasi impossible de former correctement l'agent puisqu'on met en place une probabilité de réaliser l'action voulu (malgré l'augmentation à 10 000 répétitions du jeu). En effet, même si l'agent réussit à atteindre l'objectif, la répétition de la même politique ne permettra pas forcément de réussir à nouveau puisqu'il y'aura toujours une probabilité de 3/4 de glisser (d'où la répétition non continue des victoires suivie des pertes à chaque fois). On peut également voir que les valeurs de notre table-Q sont plus faible et qu'il y a plus de valeurs nulles.

Dans le cas simple, l'environnement est déterministe, ce qui signifie que l'agent sait avec certitude ce qui se passera si il prend une action donnée. Dans ce cas, le Q-learning peut apprendre rapidement et efficacement la politique optimale.

Dans le cas complexe, l'environnement est non déterministe, ce qui signifie que l'agent ne sait pas toujours ce qui se passera si il prend une action donnée. Dans ce cas, le Q-learning peut avoir du mal à apprendre la politique optimale.

Par exemple, dans le cas de la marche, si l'action "aller à droite" a une probabilité de $\frac{2}{3}$ de mener à la falaise, le Q-learning aura du mal à apprendre que cette action est mauvaise.

Les limites du Q-learning dans les cas complexes peuvent être expliquées par les facteurs suivants :

- La complexité de l'espace d'états et d'actions : Plus l'espace d'états et d'actions est complexe, plus il est difficile pour le Q-learning d'apprendre la politique optimale.
- La non-déterminisme de l'environnement : Le Q-learning est basé sur la notion de valeur de qualité, qui est une estimation de la probabilité qu'une action mène à un résultat donné. Dans un environnement non déterministe, la valeur de qualité d'une action n'est pas toujours fiable.
- Le temps d'apprentissage : Le Q-learning peut prendre beaucoup de temps pour converger vers la politique optimale, en particulier dans les cas complexes.

Le machine learning supervisé est une autre méthode d'apprentissage automatique qui peut être utilisée pour résoudre des problèmes complexes. Dans le machine learning supervisé, l'agent est formé sur un ensemble de données de données d'entraînement. Ces données d'entraînement contiennent des exemples d'états et d'actions, ainsi que les résultats souhaités.

Dans le cas de la marche, l'agent pourrait être formé sur un ensemble de données de données d'entraînement qui contiennent des exemples de positions du marcheur, ainsi que les actions qui ont mené à la victoire.

Le machine learning supervisé est généralement plus efficace que le Q-learning dans les cas complexes. Cependant, il nécessite un ensemble de données d'entraînement de bonne qualité, ce qui peut être difficile à obtenir.