

Apprentissage automatique

apprendre à partir d'exemples

Charles Prud'homme

Charles.Prudhomme@imt-atlantique.fr

TASC (CNRS/IMT Atlantique)



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom

1 Apprentissage par renforcement

1 Apprentissage par renforcement

Apprentissage de réflexes par renforcement

Mise en situation

Vous jouez pour la première fois à un jeu dont vous ne connaissez pas les règles. Après une centaine de coups, votre adversaire annonce : “**Tu as perdu**”.

Définition

Un algorithme d'apprentissage par renforcement s'attache

- à apprendre les actions à prendre
- à partir d'expériences
- de façon à optimiser une récompense quantitative
- au cours du temps.

Apprentissage de réflexes par renforcement

Mise en situation

Vous jouez pour la première fois à un jeu dont vous ne connaissez pas les règles. Après une centaine de coups, votre adversaire annonce : “**Tu as perdu**”.

Définition

Un algorithme d'apprentissage par renforcement s'attache

- à apprendre les actions à prendre
- à partir d'expériences
- de façon à optimiser une récompense quantitative
- au cours du temps.

Suppositions

- ① L'agent ne connaît pas ou mal son environnement :
 - Ne connaît pas *a priori* quels sont les renforcements associés à chaque état ou transition
 - Ne connaît pas la topologie de l'espace
- ② L'agent ne connaît pas ou mal l'effet de ses actions dans un état donné

Suppositions

- ① L'agent ne connaît pas ou mal son environnement :
 - Ne connaît pas *a priori* quels sont les renforcements associés à chaque état ou transition
 - Ne connaît pas la topologie de l'espace
- ② L'agent ne connaît pas ou mal l'effet de ses actions dans un état donné

Défi

- Connaissance faible du monde et de ses réactions
- Mesures sur les états peuvent être imparfaites
- Renforcements pauvre en information, parfois tardif
- \Rightarrow Nécessite énormément d'interactions

Suppositions

- ① L'agent ne connaît pas ou mal son environnement :
 - Ne connaît pas *a priori* quels sont les renforcements associés à chaque état ou transition
 - Ne connaît pas la topologie de l'espace
- ② L'agent ne connaît pas ou mal l'effet de ses actions dans un état donné

Défi

- Connaissance faible du monde et de ses réactions
- Mesures sur les états peuvent être imparfaites
- Renforcements pauvre en information, parfois tardif
- \Rightarrow Nécessite énormément d'interactions

AlphGo (2016) : des 10aines de millions de parties ont été nécessaires

Suppositions

- ① L'agent ne connaît pas ou mal son environnement :
 - Ne connaît pas *a priori* quels sont les renforcements associés à chaque état ou transition
 - Ne connaît pas la topologie de l'espace
- ② L'agent ne connaît pas ou mal l'effet de ses actions dans un état donné

Défi

- Connaissance faible du monde et de ses réactions
- Mesures sur les états peuvent être imparfaites
- Renforcements pauvre en information, parfois tardif
- \Rightarrow Nécessite énormément d'interactions

Relativement inefficace mais très adaptable

Suppositions

- ① L'agent ne connaît pas ou mal son environnement :
 - Ne connaît pas *a priori* quels sont les renforcements associés à chaque état ou transition
 - Ne connaît pas la topologie de l'espace
- ② L'agent ne connaît pas ou mal l'effet de ses actions dans un état donné

Seules hypothèses valables

Le monde est

- stochastique : les actions peuvent avoir des effets non déterministes,
- stationnaire : les probabilités de transition et les renforcements restent stables

Modélisation

L'agent communique avec son environnement par trois canaux :

- Un **canal perceptif** : $s(t)$, mesure l'état dans lequel il se trouve dans l'environnement
- Un canal spécifique aux signaux de renforcement : $r(t)$, renseigne sur la qualité de l'état courant,
- Un canal qui transmet l'action de l'agent, $a(t)$, à l'environnement

Notations

Modélisation

L'agent communique avec son environnement par trois canaux :

- Un **canal perceptif** : $s(t)$, mesure l'état dans lequel il se trouve dans l'environnement
- Un canal spécifique aux signaux de renforcement : $r(t)$, renseigne sur la qualité de l'état courant,
- Un canal qui transmet l'action de l'agent, $a(t)$, à l'environnement

Notations

À l'instant t

- $s_t \in \mathcal{E}$, l'espace des états
- $r_t \in \mathcal{R}$, l'espace des signaux, $r(t) \in [-a, +b]$, $a, b \in \mathbb{R}^+$
- $a_t \in \mathcal{A}$, l'espace des actions

Modélisation

L'agent communique avec son environnement par trois canaux :

- Un **canal perceptif** : $s(t)$, mesure l'état dans lequel il se trouve dans l'environnement
- Un canal spécifique aux signaux de renforcement : $r(t)$, renseigne sur la qualité de l'état courant,
- Un canal qui transmet l'action de l'agent, $a(t)$, à l'environnement

Notations

- L' **agent** est une fonction $s_t \mapsto a_t$
- On parlera de *politique*, notée π_t
- $\pi(s, a)$: la probabilité de choisir l'action a dans l'état s .

Modélisation

L'agent communique avec son environnement par trois canaux :

- Un **canal perceptif** : $s(t)$, mesure l'état dans lequel il se trouve dans l'environnement
- Un canal spécifique aux signaux de renforcement : $r(t)$, renseigne sur la qualité de l'état courant,
- Un canal qui transmet l'action de l'agent, $a(t)$, à l'environnement

Notations

- L'**environnement** est une fonction $(s_t, a_t) \mapsto (s_{t+1}, r_t)$
- On distinguera :
 - la fonction de transition entre états, T
 - la fonction de renforcement immédiat, R

Les mesures de gain

- Pas de mesure de gain universelle
- Doit être spécifiée par problème
- En général, on choisit de cumuler le gain dans le temps
- Choix de la mesure est **déterminant**

Les mesures de gain

- **Gain cumulé avec horizon fini**

$$R_T = \sum_{t=0}^{T-1} r_t | s_0$$

- **Gain cumulé avec intérêt et horizon infini**

$$R = \sum_{t=0}^{\infty} \gamma^t r_t | s_0, \quad 0 \leq \gamma \leq 1$$

- **Gain en moyenne**

$$R_T = \frac{1}{T-1} \sum_{t=0}^{T-1} r_t | s_0$$

Les processus décisionnels de Markov

Apprentissage par renforcement

⇒ problème de *décision séquentielle dans l'incertain*.

Chaque décision de l'agent a un effet sur les décisions à suivre

⇒ les entrées ne sont pas considérées comme *i.i.d.*.

Formalisation par processus décisionnels de Markov

Généralisation de la recherche de plus court chemin dans un environnement stochastique.

La solution d'un problème de décision markovien

- n'est pas une séquence de décisions
- mais une **politique** qui spécifie l'action à prendre en chacun des états rencontrée pour maximiser l'espérance de gain.

Les processus décisionnels de Markov

Apprentissage par renforcement

⇒ problème de *décision séquentielle dans l'incertain*.

Chaque décision de l'agent a un effet sur les décisions à suivre

⇒ les entrées ne sont pas considérées comme *i.i.d.*.

Formalisation par processus décisionnels de Markov

Généralisation de la recherche de plus court chemin dans un environnement stochastique.

La solution d'un problème de décision markovien

- n'est pas une séquence de décisions
- mais une **politique** qui spécifie l'action à prendre en chacun des états rencontrée pour maximiser l'espérance de gain.

Les processus décisionnels de Markov

Apprentissage par renforcement

⇒ problème de *décision séquentielle dans l'incertain*.

Chaque décision de l'agent a un effet sur les décisions à suivre

⇒ les entrées ne sont pas considérées comme *i.i.d.*.

Formalisation par processus décisionnels de Markov

Généralisation de la recherche de plus court chemin dans un environnement stochastique.

La solution d'un problème de décision markovien

- n'est pas une séquence de décisions
- mais une **politique** qui spécifie l'action à prendre en chacun des états rencontrée pour maximiser l'espérance de gain.

Les processus décisionnels de Markov

Apprentissage par renforcement

⇒ problème de *décision séquentielle dans l'incertain*.

Chaque décision de l'agent a un effet sur les décisions à suivre

⇒ les entrées ne sont pas considérées comme *i.i.d.*.

Formalisation par processus décisionnels de Markov

Généralisation de la recherche de plus court chemin dans un environnement stochastique.

La solution d'un problème de décision markovien

- n'est pas une séquence de décisions
- mais une **politique** qui spécifie l'action à prendre en chacun des états rencontrée pour maximiser l'espérance de gain.

Les approches

Plusieurs approches pour résoudre ce problème :

Model-Based RL

- Apprendre un modèle de l'environnement (fonction de renforcement + fonction de transition)
- + \approx Apprentissage supervisé
- Ignore les interactions entre états

Les approches

Plusieurs approches pour résoudre ce problème :

Value-Based RL

- Model-based RL + *fonction d'utilité* (pour les interactions)

soit $V(s)$: espérance de gain à partir d'un **état**

soit $Q(s, a)$: espérance de gain à partir d'un **couple état-action**

- Agir sur le monde et calculer sur le long terme la qualité des états ou des couples état-action

Les approches

Plusieurs approches pour résoudre ce problème :

Policy-Based RL

Travailler sur l'espace des politiques

soit Sélection darwinienne où chaque agent correspond à une politique

soit Apprendre *directement* une politique $\pi : \mathcal{E} \mapsto \mathcal{A}$

Fonctions d'utilité

Volontés

- Optimiser la conduite sur le long terme
- sur la base de décisions locales
- ne nécessitant de recherche en avant

⇒ l'information locale doit refléter l'espérance de gain à long terme !

$$V^\pi(s) = \mathbb{E}_\pi\{R_t | s_t = s\}$$

FIGURE – Espérance de gain à partir de l'étape s en suivant la politique π .

Fonctions d'utilité

Volontés

- Optimiser la conduite sur le long terme
- sur la base de décisions locales
- ne nécessitant de recherche en avant

⇒ l'information locale doit refléter l'espérance de gain à long terme !

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

FIGURE – Espérance de gain à partir de l'étape s , en effectuant l'action a , puis en suivant la politique π .

1 Apprentissage par renforcement Quand l'environnement est connu

Préliminaires

On suppose ici connus :

- Les probabilités de transition
- Les renforcements associés
- L'agent sait ce qu'il peut atteindre dans l'environnement
- Mais ne connaît pas les fonctions d'utilités

Donc il ne connaît pas l'impact de ses décisions sur le long terme

- 1 Comment les apprendre pour une politique donnée ?
- 2 Comment approcher une politique optimale ?

Préliminaires

On suppose ici connus :

- Les probabilités de transition
- Les renforcements associés
- L'agent sait ce qu'il peut atteindre dans l'environnement
- Mais ne connaît pas les fonctions d'utilités

Donc il ne connaît pas l'impact de ses décisions sur le long terme

- ❶ Comment les apprendre pour une politique donnée ?
- ❷ Comment approcher une politique optimale ?

Préliminaires

On suppose ici connus :

- Les probabilités de transition
- Les renforcements associés
- L'agent sait ce qu'il peut atteindre dans l'environnement
- Mais ne connaît pas les fonctions d'utilités

Donc il ne connaît pas l'impact de ses décisions sur le long terme

- ① Comment les apprendre pour une politique donnée ?
- ② Comment approcher une politique optimale ?

Évaluer une politique

par propagation locale d'information

Approche simple

- Tester tous les états s
- En suivant la politique π (au moins une fois)
- Et calculer la moyenne des gains cumulés

Notation

$\pi(s, a)$: probabilité de choisir a dans l'état s alors qu'on applique la politique π

Évaluer une politique

par propagation locale d'information

Approche simple

- Tester tous les états s
- En suivant la politique π (au moins une fois)
- Et calculer la moyenne des gains cumulés

Notation

$\pi(s, a)$: probabilité de choisir a dans l'état s alors qu'on applique la politique π

Évaluation

P-ex., gain cumulé avec intérêts et horizon infini :

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{ R_t | s_t = s \} \\ &= r(s, \pi(s)) + \gamma \sum_{s'} p^\pi(s' | s_t) V^\pi(s') \end{aligned} \quad (1)$$

$$= \sum_a \pi(s, a) Q^\pi(s, a) \quad (2)$$

L'espérance de gain d'un état dépend de

- son propre renforcement
- et des espérances de gain des états qu'il peut atteindre

Évaluation

P-ex., gain cumulé avec intérêts et horizon infini :

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{R_t | s_t = s\} \\ &= r(s, \pi(s)) + \gamma \sum_{s'} p^\pi(s' | s_t) V^\pi(s') \end{aligned} \quad (1)$$

$$= \sum_a \pi(s, a) Q^\pi(s, a) \quad (2)$$

L'espérance de gain d'un état dépend de

- son propre renforcement
- et des espérances de gain des états qu'il peut atteindre

Évaluation

P-ex., gain cumulé avec intérêts et horizon infini :

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{R_t | s_t = s\} \\ &= r(s, \pi(s)) + \gamma \sum_{s'} p^\pi(s' | s_t) V^\pi(s') \end{aligned} \quad (1)$$

$$= \sum_a \pi(s, a) Q^\pi(s, a) \quad (2)$$

L'espérance de gain d'un état dépend de

- son propre renforcement
- et des espérances de gain des états qu'il peut atteindre

Évaluation

P-ex., gain cumulé avec intérêts et horizon infini :

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \{R_t | s_t = s\} \\ &= r(s, \pi(s)) + \gamma \sum_{s'} p^\pi(s' | s_t) V^\pi(s') \end{aligned} \quad (1)$$

$$= \sum_a \pi(s, a) Q^\pi(s, a) \quad (2)$$

L'espérance de gain d'un état dépend de

- son propre renforcement
- et des espérances de gain des états qu'il peut atteindre

Politique optimale

Ordre sur les politiques

$$\pi \geq \pi' \iff V^\pi(s) \geq V^{\pi'}(s), \forall s \in \mathcal{E}$$

Politique optimale π^*

Si une politique est supérieure à toutes les autres, elle est optimale et notée π^* .

Conduite optimale

Politique optimale

Ordre sur les politiques

$$\pi \geq \pi' \iff V^\pi(s) \geq V^{\pi'}(s), \forall s \in \mathcal{E}$$

Politique optimale π^*

Si une politique est supérieure à toutes les autres, elle est optimale et notée π^* .

Conduite optimale

Politique optimale

Ordre sur les politiques

$$\pi \geq \pi' \iff V^\pi(s) \geq V^{\pi'}(s), \forall s \in \mathcal{E}$$

Politique optimale π^*

Si une politique est supérieure à toutes les autres, elle est optimale et notée π^* .

Conduite optimale

Politique optimale

Ordre sur les politiques

$$\pi \geq \pi' \iff V^\pi(s) \geq V^{\pi'}(s), \forall s \in \mathcal{E}$$

Politique optimale π^*

Si une politique est supérieure à toutes les autres, elle est optimale et notée π^* .

Conduite optimale

Si l'agent dispose des $V^*(s) = \max_{\pi} V^\pi(s)$, $\forall s \in \mathcal{E}$

- pour chaque action a , faire un pas en avant
- choisir l'action avec la meilleure espérance de gain

Politique optimale

Ordre sur les politiques

$$\pi \geq \pi' \iff V^\pi(s) \geq V^{\pi'}(s), \forall s \in \mathcal{E}$$

Politique optimale π^*

Si une politique est supérieure à toutes les autres, elle est optimale et notée π^* .

Conduite optimale

Si l'agent dispose des $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$, $\forall s \in \mathcal{E}$ et $\forall a \in \mathcal{A}$

- choisir l'action avec la meilleure espérance de gain

Amélioration de politique

Théorème : Relation d'ordre sur les politiques

Soient π et π' deux politiques déterministes, telles que, pour tout état $s \in \mathcal{E}$:

$$Q^\pi(s, \pi'(s)) \geq V^\pi(s)$$

Alors la politique π' doit être au moins aussi bonne que la politique π , ce qui signifie que, pour tout état $s \in \mathcal{E}$:

$$V^{\pi'}(s) \geq V^\pi(s)$$

Un procédure de type gradient permet de choisir pour chaque état la meilleure action en fonction d'une politique

Processus itératif

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{A} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{A} \pi_2 \xrightarrow{E} \dots \xrightarrow{A} \pi^* \xrightarrow{E} V^*$$

* E : évaluation, A : amélioration.

Bilan

- + Convergence en un nombre fini d'itérations
 - Phase d'évaluation de politique est **très coûteuse**
- + On peut la limiter à un passage par état
 - En pratique, tous les états ne peuvent pas être visités...

Sans modèle du monde : la **méthode de Monte-Carlo**

→ Estimer les valeurs de probabilités de transitions et de renforcement par un échantillonnage.

Processus itératif

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{A} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{A} \pi_2 \xrightarrow{E} \dots \xrightarrow{A} \pi^* \xrightarrow{E} V^*$$

* E : évaluation, A : amélioration.

Bilan

- + Convergence en un nombre fini d'itérations
- Phase d'évaluation de politique est **très coûteuse**
- + On peut la limiter à un passage par état
- En pratique, tous les états ne peuvent pas être visités...

Sans modèle du monde : la **méthode de Monte-Carlo**

→ Estimer les valeurs de probabilités de transitions et de renforcement par un échantillonnage.

Processus itératif

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{A} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{A} \pi_2 \xrightarrow{E} \dots \xrightarrow{A} \pi^* \xrightarrow{E} V^*$$

* E : évaluation, A : amélioration.

Bilan

- + Convergence en un nombre fini d'itérations
 - Phase d'évaluation de politique est **très coûteuse**
- + On peut la limiter à un passage par état
 - En pratique, tous les états ne peuvent pas être visités...

Sans modèle du monde : la **méthode de Monte-Carlo**

→ Estimer les valeurs de probabilités de transitions et de renforcement par un échantillonnage.

Processus itératif

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{A} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{A} \pi_2 \xrightarrow{E} \dots \xrightarrow{A} \pi^* \xrightarrow{E} V^*$$

* E : évaluation, A : amélioration.

Bilan

- + Convergence en un nombre fini d'itérations
 - Phase d'évaluation de politique est **très coûteuse**
- + On peut la limiter à un passage par état
 - En pratique, tous les états ne peuvent pas être visités...

Sans modèle du monde : la **méthode de Monte-Carlo**

→ Estimer les valeurs de probabilités de transitions et de renforcement par un échantillonnage.

- 1 Apprentissage par renforcement
Évaluation de la politique par la méthode des différences temporelles

Méthode des différences temporelles

Principes

On approxime la formule : $V^\pi(s) = \mathbb{E}_\pi\{R_t | s_t = s\}$ par :

$$V(s) \leftarrow V(s) + \alpha[R_t - V(s)]$$

où R_t mesure le gain après l'instant t en partant de s , $R_t = r + V(s')$.

- $R_t - V(s)$ calcul l'erreur sur l'estimation courante = direction
- Pas besoin de connaissances *a priori* sur l'environnement
- Nécessite de ne mémoriser que $V(s)$ + calcul simple
- α est constant ou décroissant lentement

Méthode des différences temporelles

```
1: procedure TEMPORALDIFFERENCE
2:   Initialiser  $V(s)$  arbitrairement et  $\pi$  à la politique à évaluer.
3:   repeat
4:     for all épisode do
5:       Partir de  $s$ 
6:       repeat
7:         for all étape de l'épisode do
8:            $a \leftarrow$  l'action donnée par  $\pi$  pour l'état  $s$ 
9:           Exécuter  $a$ , recevoir  $r$  et  $s'$ 
10:           $V^\pi(s) \leftarrow V^\pi(s) + \alpha[r + \gamma V^\pi(s') - V^\pi(s)]$ 
11:           $s \leftarrow s'$ 
12:        until  $s$  est terminal
13:   until critère d'arrêt
```


Amélioration de politique

SARSA : méthode “sur politique”

- 1 Choisir l'action a selon une politique suivie *presque* tout le temps (procédure ε -gloutonne)
- 2 Après observation de s' et r , mettre à jour la valeur d'utilité

$$Q^\pi(s, a) \leftarrow Q^\pi(s, a) + \alpha[r + \gamma Q^\pi(s', a') - Q^\pi(s, a)]$$

Amélioration de politique

Q-learning : méthode “hors politique”

- 1 Choisir l'action a avec une procédure ε -gloutonne
- 2 Après observation de s' et r , mettre à jour la valeur d'utilité

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a)]$$

- 1 Apprentissage par renforcement
Résolution du compromis exploration contre exploitation

Exploration -vs- Exploitation

Situation

Je déménage dans une ville inconnue. J'aime manger au restaurant, mais je n'en connais aucun. Je les essaye tous une fois :

- ① je n'en favorise aucun, et continue de choisir à l'aveugle :
Exploration pure
- ② je les note tous un par un, puis je ne vais plus qu'au meilleur :
Exploitation pure

Il vaut mieux trouver un compromis entre exploration et exploitation.
Mais comment ? ⇒ résoudre un problème d'optimisation

Exploration -vs- Exploitation

Situation

Je déménage dans une ville inconnue. J'aime manger au restaurant, mais je n'en connais aucun. Je les essaye tous une fois :

- 1 je n'en favorise aucun, et continue de choisir à l'aveugle :
Exploration pure
- 2 je les note tous un par un, puis je ne vais plus qu'au meilleur :
Exploitation pure

Il vaut mieux trouver un compromis entre exploration et exploitation.
Mais comment ? \Rightarrow résoudre un problème d'optimisation

Exploration -vs- Exploitation

Situation

Je déménage dans une ville inconnue. J'aime manger au restaurant, mais je n'en connais aucun. Je les essaye tous une fois :

- 1 je n'en favorise aucun, et continue de choisir à l'aveugle :
Exploration pure
- 2 je les note tous un par un, puis je ne vais plus qu'au meilleur :
Exploitation pure

Il vaut mieux trouver un compromis entre exploration et exploitation.
Mais comment ? \Rightarrow résoudre un problème d'optimisation

Problème des bandits à bras multiples

Définition

- Il existe un ensemble de K bras, chacun défini par une distribution de récompense ν_k (dans $[0, 1]$) de loi inconnue
- À chaque pas de temps t , l'agent doit choisir un bras k_t . Il reçoit une récompense $r_t \stackrel{i.i.d}{\sim} \nu_{k_t}$
- **But** : trouver une politique de sélection des bras de manière à maximiser la somme des récompenses sur une durée donnée

Exemple de méthodes de résolution

- Méthode ε -greedy / non dirigée (*i.e.*, évalue les actions)
- Méthode basée sur la récence / dirigée (*i.e.*, + mémoire)
- **Upper Confidence Bound**

L'algorithme UCB

procédure UCB**Initialisation** : Jouer chaque bras une fois**repeat**Jouer le bras j qui maximise $\bar{x}_j + \sqrt{\frac{2 \ln n}{T_j(n)}}$ **until** fin du jeu

Où \bar{x}_j est le renforcement moyen obtenu en jouant le bras j , $T_j(n)$ le nombre de fois où le bras j a été joué et n le nombre total de tirage jusque là.

Monte-Carlo Tree Search

Que faire quand le jeu n'est pas une "répétition" mais un déplacement dans un arbre ?

Monte-Carlo Tree Search

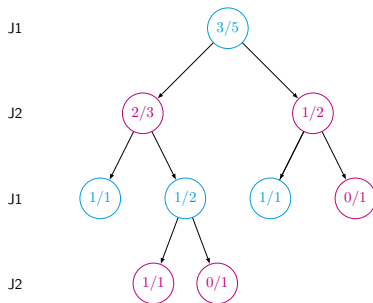
Monte-Carlo Tree Search

- Recherche arborescente, enraciné à l'état courant
- Les nœuds sont les états accessibles , les arcs sont les actions
- Une branche = séquence de coups joués
- Déterminer l'action à prendre à la racine

4 étapes

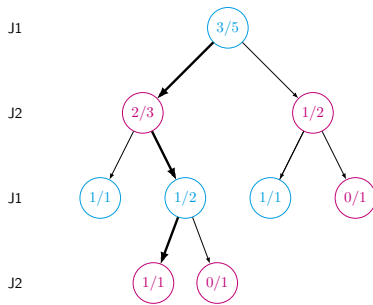
- **Sélection** : choisir le nœud le plus prometteur
- **Expansion** : créer un ou plusieurs nœuds fils et en choisir un
- **Simulation** : simuler N parties avec la politique courante pour atteindre une feuille
- **Rétro-propagation** : mise à jour des valeurs dans les nœuds

Étapes de MCTS



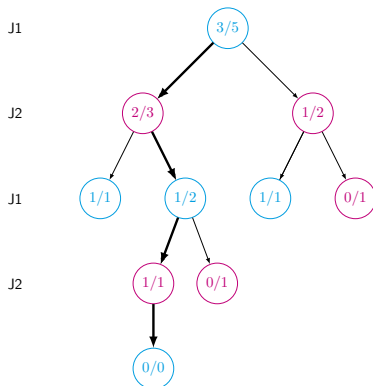
Arbre à l'étape t

Étapes de MCTS



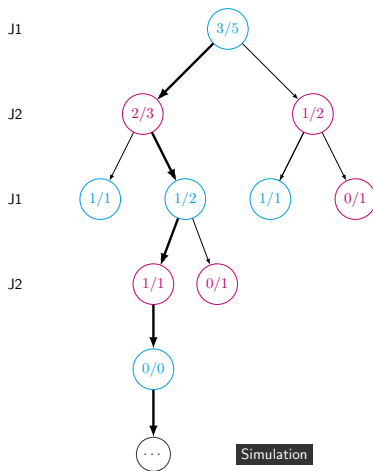
Sélection

Étapes de MCTS

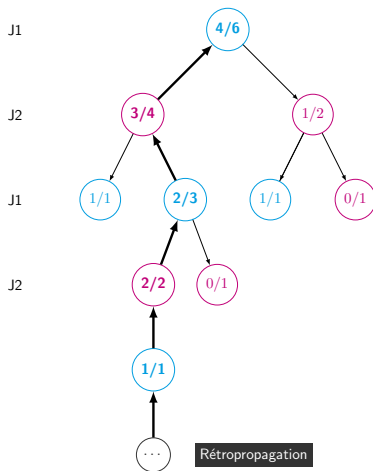


Expansion

Étapes de MCTS



Étapes de MCTS



Apprentissage automatique

apprendre à partir d'exemples

Charles Prud'homme

Charles.Prudhomme@imt-atlantique.fr

TASC (CNRS/IMT Atlantique)

