

# Shapota

## Monte Carlo Tree Search

(MCTS)

### 1. Introduction

- Toujours Apprentissage dans  
l'incertain

- Arbres de recherche dans les jeux:

+ jeux à information complète

+ jeux à deux adversaires

+ jeux à somme nulle

( $\frac{X}{X}$ , échecs, Go, ...)

→ Solution: Algorithme MinMax  
(avec l'élagage  $\alpha\beta$ )

+ efficace pour les arbres  
à faible facteur de  
branchement

+ nécessite de "bonnes"  
fonctions d'heuristiques

Exemple où la technique n'est pas du tout efficace : Go

- bn 250 en moyenne
- pas de très bonne heuristique.

## 2. Simulation de type Monte-Carlo.

1. Simuler le jeu en effectuant des mouvements aléatoires
2. Donner un score à la fin du jeu
3. Mesurer la statistique liées aux jeux gagnants
4. Jouer le mouvement avec le plus grand pourcentage de gain
5. Répéter ces actions.

Application au MinMax:  
utiliser les simulations comme  
évaluation pour les plateaux

Problèmes: - Il faut faire  
(knowement) beaucoup de  
simulations pour chaque  
plateau.

⇒ MC ont été négligées  
pendant très longtemps  
pour le jeu de Go...

- Tel quel, on va avoir tendance  
à chercher systématiquement  
de nouvelles positions alors que  
des régions (déjà explorées)  
peuvent être intéressantes.

⇒ Il faut un ~~trade off~~  
entre exploration  
et exploitation

# exploration vs exploitation

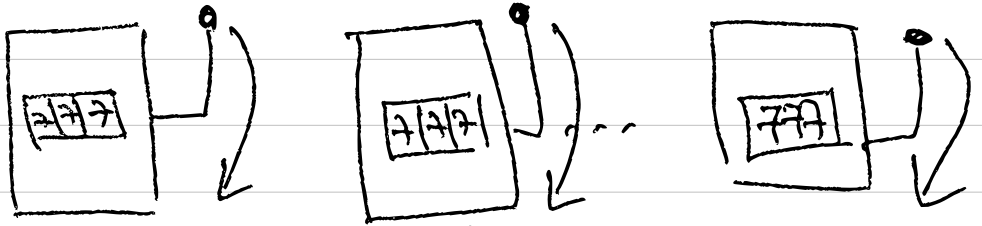
Principe: on souhaiterait utiliser le résultat de simulation déjà effectuées pour orienter la croissance de l'arbre de jeu

Exploitation: jouer le coup le plus prometteur (d'après la simulation déjà faite)

Exploration: jouer les coups pour lesquels l'incertitude est la plus haute

⇒ Théorie de bandit

# Multi Armed Bandits



1 machine = bras = action

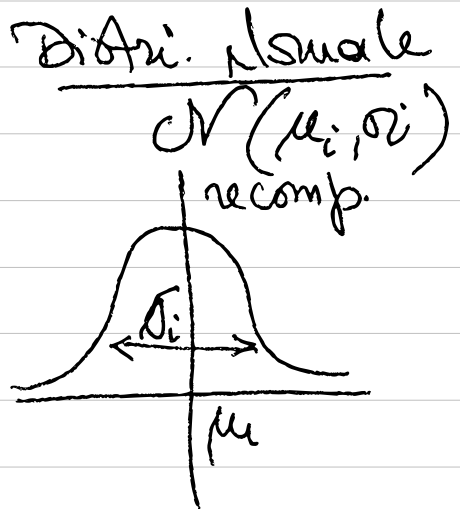
$k$  machines:  $1, 2, \dots, k$

Chaque machine  $i$  a une distribution de récompense  $\delta_i$

Bernoulli  $B(p)$

$\left\{ \begin{array}{l} 1 \text{ avec prob } p \\ 0 \text{ avec prob } 1-p \end{array} \right.$

$$\begin{aligned} \mu &= \mathbb{E}(\delta) \\ \sigma^2 &= \text{Var}(\delta) \end{aligned}$$



## Hypothèses:

- On connaît le menuisier  
les distributions  $(S_i)_{i=1, \dots, k}$

⇒ Simulation (MC):

- à chaque étape  $t$ , on choisit un bras  $i$  et on génère une récompense  $r(t) \in \mathbb{R}$  selon  $S_i$

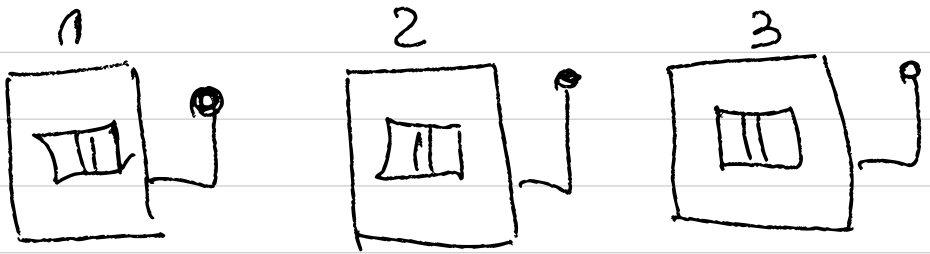
Objectifs (équivalents):

(1) Trouver le meilleur bras:  
arg max  $\mu_i$

(2) maximiser la récompense totale

$$\sum_{t \geq 1} r(t).$$

Problème: on doit faire un choix avec une information incomplète:



$$S_1 = \begin{cases} 1 & \text{prob. } .6 \\ 0 & .4 \end{cases}$$

$$S_2 = \begin{cases} 1 & \text{prob. } .55 \\ 0 & .45 \end{cases}$$

$$S_3 = \begin{cases} 1 & \text{prob. } .4 \\ 0 & .6 \end{cases}$$

$\Rightarrow$  Le meilleur choix aurait été le bras 1 mais on ne le sait pas

$\Rightarrow$  On cherche à minimiser le "regret" = la perte par rapport à la meilleure suite de bras (si on avait l'information).

$\Rightarrow$  équilibre entre exploration et exploitation  
 politique uniforme (greedy)

# The $\epsilon$ -greedy algorithm

Après  $T$  étapes

pour chaque machine, on dispose d'un échantillon

$\rightarrow \hat{\mu}_i(T)$  espérance empirique

Greedy: on joue toujours  $\hat{\mu}_i(T)$  (exploitation)

$\epsilon$ -greedy: on fixe  $\epsilon \in [0, 1]$

- avec prob.  $\epsilon$ , on choisit uniformement  
une action (exploration)

avec prob.  $1 - \epsilon$ , on choisit  
argmax  $\hat{\mu}_i(T)$  (exploitation)

Analyse des regrets: une reformulation

Regret  $(T) = R(T) = \text{diff. entre le}$   
meilleur a posteriori et ce qu'on  
obtient



$$\Rightarrow R(T) = T \cdot \mu_k - \sum_{t=1}^T r(t)$$

( $\mu_i$  est la vraie espérance et  $\hat{\mu}_i(T)$  son estimation à l'instant  $T$ )

$$\Rightarrow \min R(T) \equiv \max \sum_{t=1}^T r(t)$$

~~x~~ → ~~x~~

## UPPER Confidence Bounds (UCB)

$r(t)$  = recomp. à l'instant  $t$

$$r(i, t) = \begin{cases} r(t) & \text{si } i \text{ est choisi à } t \\ 0 & \text{sinon} \end{cases}$$

$$R(t) = T \mu_k - \sum_{\tau=1}^T r(t)$$

$$\hat{\mu}_i(T) = \frac{1}{n(i, T)} \sum_{t=1}^T r(i, t)$$

$$\boxed{\text{UCB} = \arg \max_i \hat{\mu}_i(T) + c(i, T)}$$

avec  $c(i, T) = \sqrt{\log(T) / n(i, T)}$

## Intuitions:

- si  $m(i, T)$  est petit (peu d'information)  
alors  $c(i, T)$  est grand  
 $\Rightarrow$  Il faut explorer

- si  $m(i, T)$  est grand alors  $c(i, T)$   
est petit, sauf quand  $T$  est  
très grand

Question: pourquoi  $c(i, T) = \sqrt{\frac{\log(T)}{m(i, T)}}$ .

Borne de Chernoff-Hoeffding :

$Y_1, Y_2, \dots, Y_n$  iid de m<sup>e</sup> loi que  $Y$   
- loi de grands nombres

$$\frac{1}{n} \sum Y_i \rightarrow \mathbb{E}(Y)$$

-  $\Pr\left(\left|\frac{1}{n} \sum Y_i - \mathbb{E}(Y)\right| \geq c\right) \leq 2e^{-2c^2 n^2}$

Application:

$$\mathbb{P}(|\hat{\mu}_i(T) - \mu_i| \geq c(i; T)) \\ \leq 2 e^{-2 c(i; T)^2 m(i; T)^2}$$

avec  $c(i; T) = \sqrt{\frac{\log(T)}{m(i; T)}}$ , la borne

devient  $\frac{2}{T^2} \xrightarrow{T \rightarrow \infty} 0$

Mise à jour :

$$\hat{\mu}_i(T) = \frac{1}{m(i; T)} \sum_{t=1}^T x(i; t)$$

si on choisit le tonas  $i$ :

$$\hat{\mu}_i(T+1) = \hat{\mu}_i(T) + \frac{1}{n(i; T)} [\hat{\mu}_i(T) - x(i; T+1)]$$

new = old +  $\alpha$  (old - current)  
 $\uparrow$  step size

# MCTS.

- exploration d'un arbre de jeu:

4 étapes:

1. Selection  
(d'une feuille à  
développer)



2. Expansion



3. Simulation  
(roll out)



4. Back-propagation

On met à jour  
les valeurs de  
ancêtres de la feuille  
explorée.



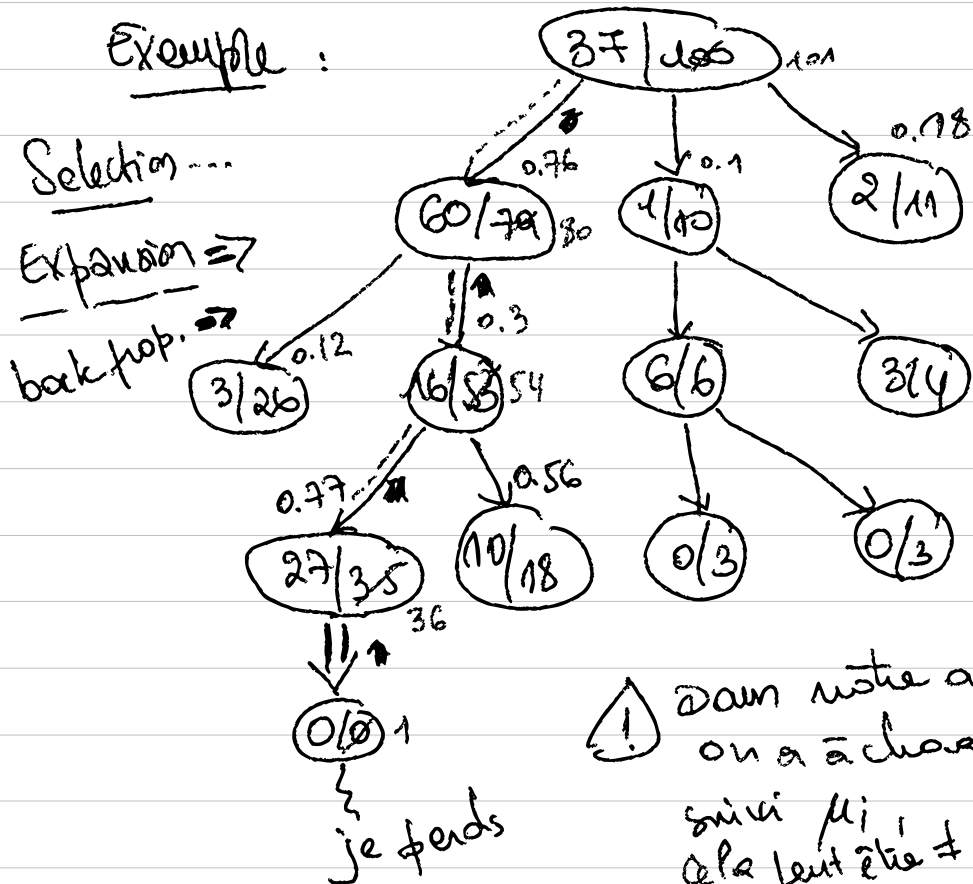
# Upper Confidence Tree = MCTS + UCB

Comment faire la Selection ?  
réponse : UCB.

$$\underline{UCB} : \underset{i}{\operatorname{argmax}} \mu_i(T) + c(i, T)$$

$$\mu_i(T) = \frac{1}{n(i, T)} \sum_{t=1}^T r(i, t)$$

Exemple :



# Simulation (roll out)



- Répéter

1.  $a_i = \text{random}(\text{actions}(s_i))$
  2.  $s_i = \text{simulate}(a_i, s_i)$
- jusqu'à  $\bar{s}_i$   $s_i$  est terminal

- retourner  $\text{valeur}(s_i)$

## Application à l'apprentissage par renforcement:

Ce qu'on veut approcher

v-values =  $v(s)$  pour chaque état  $s$

q-values =  $q(s, a)$  pour chaque paire  
(état, action) =  $(s, a)$

## Structure de données

Plays :  $S \rightarrow \mathcal{N}$

plays(s) combien de simulations contiennent l'état  $s$ ?

$$v : S \rightarrow \mathbb{R}$$

$v(s)$  = estimation de  $v^*(s)$

$$q : S, A \rightarrow \mathbb{R}$$

$q(s, a)$  = estimation de  $q^*(s, a)$

Exercice : - MCTS pour  $U$   
- MCTS pour  $q$ .

