

07/02/2013

Bandits = à la fin de chaque tour, on observe seulement le gain (ou la perte) associé à l'action qu'on a choisie.

Dès un premier temps, on va relâcher cette hypothèse en supposant qu'on a accès à tous les gains qu'on aurait pu encaisser si on avait procédé différemment ("information parfaite").

I/ Apprentissage séquentiel robuste : combinaison d'actions en information parfaite

1) Formulation du problème

À chaque tour  $t \in \mathbb{N}^*$ , le statisticien choisit une action  $I_t \in \{1, \dots, K\}$  selon le protocole suivant.

Jeu de prédiction: à chaque tour  $t \in \mathbb{N}^*$ ,

1) Le statisticien choisit  $p_t \in \Delta(K) = \mathcal{M}_1^+(\{1, \dots, K\})$  en fonction des données disponibles  $(I_s, p_s), s=1, \dots, t-1$ .  $\mapsto p_t$  révélé à l'environnement.

2) Simultanément:

- Le statisticien tire  $I_t \sim p_t$  (conditionnellement au passé)

- L'environnement choisit  $\underline{l}_t = (l_{i,t})_{1 \leq i \leq K} \in [0, M]^K$  en fonction des données disponibles  $(I_s, p_s), s=1, \dots, t-1$  et même  $p_t$  (mais pas  $I_t$ ).

3) Le statisticien encourt la perte  $l_{I_t, t}$  et observe le vecteur  $\underline{l}_t$ ; l'environnement observe  $I_t$ .

Objectif de prédiction: minimiser le "regret"

$$R_T = \sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t}$$

Plus précisément, on cherche des stratégies/politiques (= choix séquentiel des  $p_t$ ) qui assurent que

des  $p_t$ ) qui assurent que  $\left. \begin{array}{l} \bullet \mathbb{E}[R_T] \leq o(T) \quad (\text{ex: } \mathbb{E}[R_T] \leq o(\sqrt{T \ln K})) \\ \bullet \text{p.s., } \limsup_{T \rightarrow +\infty} \frac{R_T}{T} \leq 0. \end{array} \right\} \begin{array}{l} \text{pour toute suite} \\ \text{de fonctions } l_t(\cdot) \\ \text{(POUR TOUT ADVERSAIRE)} \end{array}$

2) Le prédicteur par pondération exponentielle

Rappel: pour espérer vérifier  $\mathbb{E}[R_T] \leq o(T)$ , il faut convexifier!

Autrement dit:  $p_t = \text{diac}$  est prosrit.

# Algorithme : EWA( $\eta$ ) ["Exponentially Weighted Average"]

• Paramètre :  $\eta > 0$

• et chaque date  $t \in \mathbb{N}^*$ ,

$$P_{i,t} = \frac{\exp(-\eta \sum_{s=1}^{t-1} l_{i,s})}{\sum_{j=1}^K \exp(-\eta \sum_{s=1}^{t-1} l_{j,s})}, \quad 1 \leq i \leq K.$$

Lem :  $P_1 = (\frac{1}{K}, \dots, \frac{1}{K})$

celui ne dépend que des pertes passées  $l_s, s=1, \dots, t-1$ .

$\eta$  à calibrer

## Proposition 1: borne de regret de EWA( $\eta$ )

Supposons que les pertes  $l_{i,t}$  sont à valeurs  $[0, M]$ .

Alors, p.s.,

$$\begin{aligned} \sum_{t=1}^T \sum_{i=1}^K P_{i,t} l_{i,t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} &\leq \frac{\ln K}{\eta} + \frac{\eta T M^2}{8} \\ &\leq \sqrt{\frac{T}{2} \ln K} \quad \text{pour le choix} \\ &\quad \text{de } \eta = \frac{1}{M} \sqrt{\frac{8 \ln K}{T}}. \end{aligned}$$

Preuve : on note  $P_t \cdot l_t = \sum_{i=1}^K P_{i,t} l_{i,t}$

Pour tout  $t \in \{1, \dots, T\}$ , on a :

$$\begin{aligned} P_t \cdot l_t &= -\frac{1}{\eta} \ln \left( \sum_{i=1}^K P_{i,t} e^{-\eta l_{i,t}} \right) + \frac{1}{\eta} \ln \left( \sum_{i=1}^K P_{i,t} e^{-\eta (l_{i,t} - P_t \cdot l_t)} \right) \\ &\stackrel{\text{def de } P_t}{\leq} -\frac{1}{\eta} \ln \left( \frac{\sum_{i=1}^K \exp(-\eta \sum_{s=1}^t l_{i,s})}{\sum_{i=1}^K \exp(-\eta \sum_{s=1}^{t-1} l_{i,s})} \right) + \frac{\eta M^2}{8} \leq \frac{\eta^2 M^2}{8} \quad \text{d'après le lemme de Hoeffding} \end{aligned}$$

En sommant sur  $t=1, \dots, T$  et en posant  $W_t = \frac{1}{K} \sum_{i=1}^K \exp(-\eta \sum_{s=1}^{t-1} l_{i,s})$  ( $W_1 = 1$ )

$$\begin{aligned} \sum_{t=1}^T P_t \cdot l_t &\leq -\frac{1}{\eta} \sum_{t=1}^T \ln \left( \frac{W_{t+1}}{W_t} \right) + \frac{\eta T M^2}{8} \\ &= -\frac{1}{\eta} \ln \frac{W_{T+1}}{W_1} + \frac{\eta T M^2}{8} \\ &= -\frac{1}{\eta} \ln \left( \frac{1}{K} \sum_{i=1}^K \exp(-\eta \sum_{t=1}^T l_{i,t}) \right) + \frac{\eta T M^2}{8} \end{aligned}$$

$$\leq -\frac{1}{\gamma} \ln \left( \max_{1 \leq i \leq K} \exp \left( -\gamma \sum_{t=1}^T l_{i,t} \right) \right) + \frac{\ln K}{\gamma} + \frac{\gamma T M^2}{8}$$

$$= \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} + \frac{\ln K}{\gamma} + \frac{\gamma T M^2}{8}$$

cf explication bayésienne + formule de dualité sur la Kullback par avantage d'intuition.

↳ Rappel Hoeffding :

Soit  $Z$  une v.a. centrée tq  $a \leq Z \leq b$  p.s.

Alors:  $\forall \lambda \in \mathbb{R}, \ln \mathbb{E}[e^{\lambda Z}] \leq \frac{\lambda^2 (b-a)^2}{4}$ . ("Z sous-gaussien de facteur de variance  $\frac{(b-a)^2}{4}$ ")

Conséquences de la proposition :

• Majoration du regret en espérance : comme  $\mathbb{E}[l_{I_t, t}] = \mathbb{E}[\mathbb{E}[l_{I_t, t} | \mathcal{F}_{t-1}]]$   
on a:  $\mathbb{E}\left[\sum_{t=1}^T l_{I_t, t}\right] = \mathbb{E}\left[\sum_{t=1}^T p_t \cdot l_t\right] = \mathbb{E}[p_t \cdot l_t]$

$$\text{d'où } \mathbb{E}[R_T] \leq \sqrt{\frac{T \ln K}{2}} = o(T).$$

↑  
ergodique  
par  $(I_1, \dots, I_{t-1})$

• Majoration p.s.

concentration de martingales (Hoeffding - Azuma)

$$\forall T \geq 1, \text{ avec proba } \geq 1 - \delta, \sum_{t=1}^T l_{I_t, t} \leq \sum_{t=1}^T p_t \cdot l_t + M \sqrt{\frac{T \ln \frac{1}{\delta}}{2}}$$

En posant  $\delta_T = \frac{6/\pi^2}{T^2}$  et en appliquant Borel-Cantelli, il vient :

$$\text{p.s., pour } T \text{ suffisamment grand, } R_T \leq \sqrt{\frac{T \ln K}{2}} + M \sqrt{\frac{T \ln \frac{\pi^2 T^2}{6}}{2}}$$

$$\text{d'où p.s. } \lim_{T \rightarrow \infty} \frac{R_T}{T} \leq 0.$$

N.B. On peut faire mieux :  
 $\sqrt{T \ln \ln T}$  (itéré)

• Prédiction avec avis d'experts et perte convexe

Si  $l_{i,t} = l(a_{i,t}, y_t)$  avec  $l(\cdot, y)$  convexe  $\forall y \in \mathcal{Y}$ , alors

$$\sum_{t=1}^T l\left(\sum_{i=1}^K p_{i,t} a_{i,t}, y_t\right) \leq \sum_{t=1}^T p_{i_t} \cdot l_t \leq \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} + \sqrt{\frac{T \ln K}{2}} \text{ p.s.}$$

↑  
autre moyen de convexifier !  
(plus besoin de randomiser)

3) Calibration séquentielle du paramètre de température

Pb : le choix de  $\gamma = \frac{1}{M} \sqrt{\frac{8 \ln K}{T}}$  dépend de  $T$  (et de  $M$ ).

On va choisir  $\gamma_t$  (éventuellement en fonction des données observées).

Algo: 
$$P_{i,t} = \frac{\exp(-\eta_t \sum_{s=1}^{t-1} l_{i,s})}{\sum_{j=1}^K \exp(-\eta_t \sum_{s=1}^{t-1} l_{j,s})}, \quad 1 \leq i \leq K.$$

Proposition 2: p.s., si la suite  $(\eta_t)_{t \geq 1}$  est décroissante, alors :

$$\sum_{t=1}^T P_t \cdot l_t - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} \leq \frac{\ln K}{\eta_{T+1}} + \underbrace{\sum_{t=1}^T \frac{1}{\eta_t} \ln \left( \sum_{i=1}^K P_{i,t} e^{-\eta_t (l_{i,t} - P_t \cdot l_t)} \right)}_{=: \delta_t}$$

Preuve: adaptation de la preuve ↑

$$P_t \cdot l_t \stackrel{\text{def } P_t}{=} -\frac{1}{\eta_t} \ln \left( \frac{\sum_{i=1}^K \exp(-\eta_t \sum_{s=1}^t l_{i,s})}{\sum_{i=1}^K \exp(-\eta_t \sum_{s=1}^{t-1} l_{i,s})} \right) + \delta_t$$

$$= -\frac{1}{\eta_t} \ln \left( \frac{W'_{t+1}}{W_t} \right) + \delta_t \quad \text{où } W'_{t+1} = \frac{1}{K} \sum_{i=1}^K \exp \left( -\eta_t \sum_{s=1}^t l_{i,s} \right)$$

$$= \frac{\ln W_t}{\eta_t} - \frac{\ln W_{t+1}}{\eta_{t+1}} + \left( \frac{\ln W_{t+1}}{\eta_{t+1}} - \frac{\ln W'_{t+1}}{\eta_t} \right)$$

↑ au lieu de  $\eta_{t+1} \leq \eta_t$   
 $\leq 0$  d'après Jensen et car  $\eta_{t+1} \leq \eta_t$

En sommant sur  $t=1, \dots, T$ , il vient :

$$\sum_{t=1}^T P_t \cdot l_t \leq -\frac{1}{\eta_{T+1}} \ln W_{T+1} + \sum_{t=1}^T \delta_t \quad (\text{car } \ln W_1 = \ln 1 = 0)$$

$$\leq \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} + \frac{\ln K}{\eta_{T+1}} + \sum_{t=1}^T \delta_t \quad \blacksquare$$

Conséquences :

- $\eta_t = \frac{c}{M} \sqrt{\frac{\ln K}{t}}$  et majoration Hoeffding  $\Rightarrow$  regret conditionnel  $\leq \square \sqrt{T \ln K}$   
 $\forall T \geq 1$
- Autres choix plus fins de  $\eta_t$ , cf par ex. Cesa-Bianchi, Mannor et Eloty (2007) par une borne en variance : regret conditionnel  $\leq \square \sqrt{\sum_{t=1}^T V_t \ln K} + \dots$
- Ex. de majoration plus fine que Hoeffding : majoration de type Bennett

Rappel : Si  $Z$  est une v.a. tq  $Z \leq b$  p.s. ( $b > 0$ ) et  $\mathbb{E}[Z^2] = \nu$ ,

alors :  $\forall \lambda > 0, \ln \mathbb{E}[e^{\lambda(Z - \mathbb{E}Z)}] \leq \frac{\nu}{b^2} \phi(\lambda b)$  où  $\phi(x) = e^x - x - 1$ .

Le contrôle sous-poisson sera utile ensuite. Mentionnons une conséquence:

Si  $Z \leq 0$  p.s. et  $\mathbb{E}[Z^2] = \nu$ , alors

$$(*) \quad \forall \lambda > 0, \text{ on } \mathbb{E}[e^{\lambda(Z - \mathbb{E}Z)}] \leq \frac{\lambda^2 \nu}{2} \quad (\text{passage à la limite } \lambda \rightarrow 0).$$

Application: majoration de  $\delta_t$  de type Bennett

$$\delta_t = \frac{1}{\eta_t} \ln \left( \sum_{i=1}^K p_{i,t} e^{-\eta_t (l_{i,t} - p_t \cdot l_t)} \right) = \frac{1}{\eta_t} \ln \mathbb{E}_{i \sim p_t} \left[ e^{-\eta_t (l_{i,t} - \mathbb{E}_{i \sim p_t} l_{i,t})} \right]$$

d'après (\*)  
avec  $\begin{cases} Z = l_{i,t} \\ i \sim p_t \end{cases} \rightarrow \leq \frac{\eta_t}{2} \sum_{i=1}^K p_{i,t} l_{i,t}^2$  car  $\underline{l_{i,t} \geq 0}$ .

## II / Information imparfaite : les bandits antagonistes

### 1) Formulation du problème

⚠️ Désormais, le statisticien n'observe plus  $l_{i,t}$  pour  $i \neq I_t$ .

Protocole de décision (bandits antagonistes à  $K$  bras) : pour chaque  $t \in \mathbb{N}^+$ ,

- 1) Le statisticien choisit et révèle  $p_t \in \mathcal{M}_1^+(\{1, \dots, K\})$  en fonction des données disponibles  $(l_{I_s, s}, I_s)_{s \leq t-1}$ .
- 2) Simultanément:
  - Le statisticien tire  $I_t \sim p_t$  (conditionnellement au passé).
  - L'environnement choisit  $\underline{l}_t = (l_{i,t})_{1 \leq i \leq K} \in [0, M]^K$  en fonction des données disponibles  $(I_s, p_s)_{s \leq t-1}$  et même  $p_t$  (mais pas  $I_t$ ).
- 3) Le statisticien encout et observe la perte  $l_{I_t, t}$  (les  $l_{j,t}$ ,  $j \neq I_t$ , restent cachés); l'environnement observe  $I_t$ .

Objectifs: minimiser le regret  $R_T = \sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t}$

Plus précisément, on cherche des stratégies (choix séquentiels des  $p_t$ ) telles que, pour toute suite de fonctions  $l_t(\cdot)$  (i.e. pour tout adversaire),

- $\mathbb{E}[R_T] \leq o(T)$  (on obtiendra  $\mathbb{E}[R_T] \leq O(\sqrt{TK \ln K})$ )
- p.s.,  $\lim_{T \rightarrow +\infty} \frac{R_T}{T} \leq 0$  (on obtiendra  $\overline{\lim}_{T \rightarrow +\infty} \frac{R_T}{\sqrt{TK \ln T}} \leq 0$ )

Dans un premier temps, on va s'intéresser à une quantité plus faible:

$$\mathbb{E} \left[ \sum_{t=1}^T l_{I_t, t} \right] - \min_{1 \leq i \leq K} \mathbb{E} \left[ \sum_{t=1}^T l_{i,t} \right] \leftarrow \text{le "pseudo-regret"}$$