

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

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

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

Application: majoration de δ_t de type Bennett

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

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

II / Information imparfaite: les bandits antagonistes

1) Formulation du problème

⚠ Adversaire, 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 $\underline{l}_t = (l_{i,t})_{1 \leq i \leq K} \in [0, 1]^K$ en fonction des données disponibles $(I_s, p_s)_{s \leq t-1}$ et même p_t (mais pas I_t).
- 3) Le statisticien encourt et observe la perte $l_{I_t, t}$ (les $l_{j,t}, j \neq I_t$, restent cachés); l'environnement observe I_t .

Objectif: 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., $\overline{\lim}_{T \rightarrow +\infty} \frac{R_T}{T} \leq 0$ (on obtiendra $\overline{\lim}_{T \rightarrow +\infty} \frac{R_T}{\sqrt{T \ln \ln T}} \leq 0$)

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

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

2) Estimation des pertes et algorithme associé

On estime \underline{l}_t par $\tilde{\underline{l}}_t$ défini par $\tilde{l}_{i,t} = \frac{l_{I_t,t}}{P_{i,t}} \mathbb{1}_{\{I_t=i\}}$

↑ Ne dépend que des données observées jusqu'à la fin du tour t .

On a $\tilde{l}_{i,t} = \frac{l_{i,t}}{P_{i,t}} \mathbb{1}_{\{I_t=i\}}$ d'où $E[\tilde{l}_{i,t} | \mathcal{F}_{t-1}] = \frac{l_{i,t}}{P_{i,t}} \underbrace{P(I_t=i | \mathcal{F}_{t-1})}_{= P_{i,t}} = l_{i,t}$

Algorithme:

$$P_{i,t} = \frac{\exp(-\eta_t \sum_{s=1}^{t-1} \tilde{l}_{i,s})}{\sum_{j=1}^K \exp(-\eta_t \sum_{s=1}^{t-1} \tilde{l}_{j,s})}, \quad 1 \leq i \leq K \quad (\text{N.B. } P_1 = (\frac{1}{K}, \dots, \frac{1}{K}))$$

3) Majoration du pseudo-regret

Théorème: Calibré avec $\eta_t = \sqrt{\frac{\ln K}{tK}}$, l'algorithme Exp3 vérifie

$$E\left[\sum_{t=1}^T l_{I_t,t}\right] - \min_{1 \leq i \leq K} \left[\sum_{t=1}^T l_{i,t}\right] \leq (1+\sqrt{2}) \cdot \sqrt{TK \ln K} \quad \left[\text{N.B. pertes } l_{i,t} \text{ dans } [0,1] \right]$$

Preuve: d'après prop 2 et la majoration de S_t de type Bennett (page 4), on a:

$$(*) \quad \sum_{t=1}^T P_t \cdot \tilde{\underline{l}}_t \leq \min_{1 \leq i \leq K} \sum_{t=1}^T \tilde{l}_{i,t} + \frac{\ln K}{\eta_{T+1}} + \sum_{t=1}^T \frac{\eta_t}{c} \sum_{i=1}^K P_{i,t} \tilde{l}_{i,t}^2 \quad \text{p.s.} \\ = \sum_{i=1}^K \frac{l_{i,t}^2}{P_{i,t}}$$

En prenant l'espérance à gauche et à droite de l'inégalité, et en remarquant que $E[\tilde{l}_{i,t}] = E[E[\tilde{l}_{i,t} | \mathcal{F}_{t-1}]] = E[l_{i,t}]$ ainsi que $E[P_t \cdot \tilde{\underline{l}}_t] = E[P_t \cdot \underline{l}_t]$, il vient:

$$E\left[\sum_{t=1}^T P_t \cdot \underline{l}_t\right] \stackrel{\text{Jensen}}{\leq} \min_{1 \leq i \leq K} E\left[\sum_{t=1}^T l_{i,t}\right] + \frac{\ln K}{\eta_{T+1}} + \sum_{t=1}^T \frac{\eta_t}{c} E\left[\frac{l_{I_t,t}^2}{P_{I_t,t}}\right]$$

En remarquant que $E\left[\frac{l_{I_t,t}^2}{P_{I_t,t}}\right] \stackrel{l_{i,t} \in [0,1]}{\leq} E\left[\frac{1}{P_{I_t,t}}\right] = E\left[\sum_{i=1}^K P_{i,t} \frac{1}{P_{i,t}}\right] = K$

et en utilisant le choix $\eta_t = \sqrt{\frac{\ln K}{tK}}$, on obtient finalement:

$$\mathbb{E} \left[\sum_{t=1}^T l_{I_{t,t}} \right] - \min_{1 \leq i \leq K} \mathbb{E} \left[\sum_{t=1}^T l_{i,t} \right] \leq \sqrt{(T+1)K \ln K} + \frac{K}{2} \sum_{t=1}^T \sqrt{\frac{\ln K}{tK}}$$

$$\leq (1+\sqrt{2}) \sqrt{TK \ln K} \quad \text{car} \quad \sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T} \quad \blacksquare$$

Remarque sur les vitesses :

- Information parfaite : $\mathbb{E}[R_T] \leq \sqrt{\frac{T}{2} \ln K}$
- Bandits : pseudo-regret $\leq (1+\sqrt{2}) \sqrt{TK \ln K}$

Même vitesse en T (\sqrt{T}) mais complexité supérieure en K pour les bandits.

C'est naturel : en divisant par T , on a les majorations suivantes sur le (pseudo) regret :

- Info parfaite : $\square \sqrt{\frac{\ln K}{T}}$
- Bandits : $\square \sqrt{\frac{\ln K}{(T/K)}}$ ← K fois moins d'observations que dans le cas "info parfaite".

Autre remarque : pseudo-regret = $\mathbb{E}[\text{regret}]$ pour un adversaire aveugle

- Def : un adversaire $(l_t(\cdot))_{t \geq 1}$ est dit "aveugle" ("oblivious") si les fonctions l_t sont constantes : $\forall t \geq 1, \forall a, b, \underbrace{l_t(a) = l_t(b)}_{\text{pertes } l_{i,t} \text{ déterministes}}$
(l'adversaire ne se sert pas du passé)

- Pour un adversaire aveugle, on a :

$$\underbrace{\mathbb{E} \left[\sum_{t=1}^T l_{I_{t,t}} \right] - \min_{1 \leq i \leq K} \mathbb{E} \left[\sum_{t=1}^T l_{i,t} \right]}_{\text{pseudo-regret}} = \mathbb{E} \left[\underbrace{\sum_{t=1}^T l_{I_{t,t}} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t}}_{\text{regret}} \right]$$

↑
déterministe

Le thm précédent fournit donc aussi une majoration du regret pour des adversaires aveugles, mais pas pour des adversaires antagonistes.

4) Amélioration : majoration du regret avec grande proba

On donne seulement quelques pistes.

On considère un adversaire quelconque (potentiellement antagoniste).

En suite obtenir des bornes de la forme :

$$\sum_{t=1}^T l_{I_t, t} \leq \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} + \Delta_{T, K}(\delta) \text{ avec probabilité } \geq 1 - \delta.$$

Puisque le choix précédent de $\tilde{l}_{i, t} = \frac{l_{i, t}}{p_{i, t}} \mathbb{1}_{\{I_t=i\}}$ conduit à $p_t \cdot \tilde{l} = \sum_{i=1}^K p_{i, t} \tilde{l}_{i, t} = l_{I_t, t}$, on pourrait vouloir partir de l'inégalité (*) page 6 et la combiner avec une majoration de la forme : $\min_{1 \leq i \leq K} \sum_{t=1}^T \tilde{l}_{i, t} \leq \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} + \dots$ avec gde proba

Problème : à chaque date t , la variance $\text{Var}(\tilde{l}_{i, t} | \mathcal{F}_{t-1}) = \frac{l_{i, t}^2 (1 - p_{i, t})}{p_{i, t}}$ peut être grande (de l'ordre de $1/p_{i, t}$). D'ici plusieurs idées :

Solution n°1 (1^{ère} partie du papier de Auer, Cesa-Bianchi, Freund et Leobuscher 2002)

On mélange l'EWA avec l'uniforme :

$$p_{i, t} = (1 - \gamma) \frac{e^{-\gamma \tilde{L}_{i, t-1}}}{\sum_{j=1}^K e^{-\gamma \tilde{L}_{j, t-1}}} + \frac{\gamma}{K}, \quad 1 \leq i \leq K \quad \left. \begin{array}{l} \text{algorithme Exp3} \\ \text{"Exponential Weights} \\ \text{for Exploration} \\ \text{and Exploitation"} \end{array} \right\}$$

où $\tilde{L}_{i, t-1} = \sum_{s=1}^{t-1} \tilde{l}_{i, s}$ et $\gamma > 0$ paramètre à calibrer.

Malheureusement, les meilleurs choix de γ ne donnent qu'une vitesse

$$\Delta_{T, K}(\delta) \approx T^{2/3} \text{ sous-optimale (pour un contrôle en grande proba).}$$

Solution n°2 (2^e partie du papier)

Algorithme Exp3.P : ("P" = "probability")

• Paramètres : $\eta > 0$ et $\beta, \gamma \in]0, 1[$.

• Initialisation : $p_1 = (\frac{1}{K}, \dots, \frac{1}{K})$

• À chaque date $t \geq 1$,

□ Tirer aléatoirement $I_t \sim p_t$ conditionnellement à $\mathcal{F}_{t-1} = \sigma(I_1, \dots, I_{t-1})$.

□ Observer $l_{I_t, t}$ et calculer $\tilde{l}_{i, t} = \frac{(1 - l_{i, t}) \mathbb{1}_{\{I_t=i\}} + \beta}{p_{i, t}}$, $i=1, \dots, K$

□ Choisir le vecteur de poids p_t selon :

$$p_{i, t} = (1 - \gamma) \frac{e^{-\eta \tilde{L}_{i, t-1}}}{\sum_{j=1}^K e^{-\eta \tilde{L}_{j, t-1}}} + \frac{\gamma}{K}, \quad 1 \leq i \leq K.$$

Remarques : • Les pertes ont été transformées en gains : $l_{i, t} \in [0, 1] \rightarrow g_{i, t} = 1 - l_{i, t} \in [0, 1]$
C'est utile pour contrôler les queues de distribution à droite des $\sum_{t=1}^T \tilde{l}_{i, t}$ (car $\tilde{l}_{i, t} \leq 1$)

- On choisit volontairement les estimateurs des $l_{i,t}$:

$$E[\tilde{l}_{i,t} | \mathcal{F}_{t-1}] = l_{i,t} - \frac{\beta}{P_{i,t}}$$

Cela permet de majorer uniformément en $i \in \{1, \dots, K\}$ les $\sum_{t=1}^T \tilde{l}_{i,t}$:

$$\forall i, \sum_{t=1}^T \tilde{l}_{i,t} \leq \sum_{t=1}^T l_{i,t} + \frac{1}{\beta} \ln\left(\frac{K}{\delta}\right), \text{ avec proba } \geq 1-\delta$$

- On utilise toujours le mélange avec l'uniforme pour éviter que les pertes estimées $\tilde{l}_{i,t}$ soient trop négatives (utile pour contrôler les \log -différence

$$\frac{1}{\eta} \ln E_{i \text{ i.i.d. } \omega_t} \left[e^{-\eta (\tilde{l}_{i,t} - \omega_t \cdot \tilde{l}_t)} \right]$$

où $\omega_t = \frac{P_t - \delta(\frac{1}{K}, \dots, \frac{1}{K})}{1-\delta} = \text{EWA}$

Boone: pour des choix de la forme (avec c_1, c_2, \dots, c_5 constantes bien choisies)

$$\eta = c_1 \sqrt{\frac{\ln K}{TK}}, \quad \beta = c_2 \sqrt{\frac{\ln K}{TK}} \quad \text{et} \quad \gamma = c_3 \sqrt{\frac{K \ln K}{T}}$$

$$\sum_{t=1}^T l_{i,t} \leq \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} + c_4 \sqrt{\frac{TK}{\ln K}} \ln\left(\frac{1}{\delta}\right) + c_5 \sqrt{TK \ln K}$$

avec proba $\geq 1-\delta$

Cf. l'article de survol de Bubeck et Cesa-Bianchi (2012) pour une preuve.

III - Bornes inférieures

1) Borne inférieure en information parfaite: $\sqrt{\frac{T}{2}} \ln K$ (pertes bornées $0 \leq l_{i,t} \leq 1$)

On va montrer que la borne supérieure $\sqrt{\frac{T}{2}} \ln K$ obtenue sur le regret de l'EWA en information parfaite est minimax-optimale quand $T, K \rightarrow +\infty$ au sens suivant:

Proposition: On pose $R_{K,T}^{\text{minimax}} = \inf_{(P_t(\cdot))_{t \geq 1}} \sup_{(l_t(\cdot))_{t \geq 1}} E \left[\sum_{t=1}^T l_{i,t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} \right]$

Alors:

$$\lim_{K \rightarrow +\infty} \lim_{T \rightarrow +\infty} \frac{R_{K,T}^{\text{minimax}}}{\sqrt{\frac{T}{2}} \ln K} \geq 1.$$

N.B.: L'inf est pris sur toutes les stratégies (i.e. choix seq. des P_t en fonction du passé). Le sup s'étend sur tous les adversaires (i.e. choix seq. des l_t) tels que $0 \leq l_{i,t} \leq 1$.

Preuve (partielle) :

On se ramène d'abord à des pertes déterministes puis stochastiques i.i.d. :

$$\inf_{(P_t(\cdot))_{t \geq 1}} \sup_{(\underline{l}_t(\cdot))_{t \geq 1}} \mathbb{E}_{P_A} [\text{regret}(T)] \geq \inf_{(P_t(\cdot))_{t \geq 1}} \sup_{\underline{l}_1, \dots, \underline{l}_T \in [0,1]^K} \mathbb{E}_{P_A} [\text{regret}(T)]$$

↑ espérance par rapport à la randomisation de l'algorithme.
↑ suites déterministes (\underline{l}_t ne dépend pas des observations passées)

$$\geq \inf_{(P_t(\cdot))_{t \geq 1}} \int_{[0,1]^{KT}} \mathbb{E}_{P_A} [\text{regret}(T)] dQ(\underline{l}_1, \dots, \underline{l}_T)$$

pour toute proba Q sur $[0,1]^{KT}$.

$$= \inf_{(P_t(\cdot))_{t \geq 1}} \mathbb{E}_{P_A \otimes Q} [\text{regret}(T)] \text{ par Fubini.}$$

On choisit Q telle que $(\underline{l}_{i,t})_{\substack{1 \leq i \leq K \\ 1 \leq t \leq T}} \stackrel{\text{iid}}{\sim} \mathcal{B}(1/2)$ (i.e. $Q = \mathcal{B}(1/2)^{\otimes TK}$).

Faisons une stratégie $(P_t(\cdot))_{t \geq 1}$

et nous avons :

$$\begin{aligned} \mathbb{E}_{P_A \otimes Q} [\text{regret}(T)] &= \mathbb{E}_{P_A \otimes Q} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] \\ &= \frac{T}{2} - \mathbb{E}_{P_A \otimes Q} \left[\min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] \quad \text{car } \mathbb{E}[l_{I_t, t}] = \frac{1}{2} \\ &= \frac{1}{2} \mathbb{E}_Q \left[\max_{1 \leq i \leq K} \sum_{t=1}^T \underbrace{(1 - 2l_{i, t})}_{=: E_{i, t}} \right] \end{aligned}$$

N.B. $E_{i, t}$ Rademacher iid ($E_{i, t} = \pm 1$ avec proba $1/2$).

$$\text{Donc } \frac{\mathbb{E}_{P_A \otimes Q} [\text{regret}(T)]}{\sqrt{\frac{T}{2} \ln K}} = \frac{1}{\sqrt{2 \ln K}} \mathbb{E}_Q \left[\max_{1 \leq i \leq K} \frac{1}{\sqrt{T}} \sum_{t=1}^T E_{i, t} \right]$$

$$\xrightarrow{T \rightarrow +\infty} \mathbb{E} \left[\max_{1 \leq i \leq K} Z_i \right] \text{ avec } Z_1, \dots, Z_K \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$$

$$\text{car (eventuellement) } \frac{1}{\sqrt{T}} \sum_{t=1}^T E_{i, t} \xrightarrow[T \rightarrow +\infty]{\mathcal{L}} (Z_1, \dots, Z_K)$$

$$\begin{aligned} \text{Donc } \lim_{K \rightarrow +\infty} \lim_{T \rightarrow +\infty} \frac{\mathbb{E}_{P_A \otimes Q} [\text{regret}(T)]}{\sqrt{\frac{T}{2} \ln K}} &\geq \lim_{K \rightarrow +\infty} \frac{\mathbb{E} \left[\max_{1 \leq i \leq K} Z_i \right]}{\sqrt{2 \ln K}} \\ &= 1. \end{aligned}$$

D'où le résultat. ■

Rem 1: Contre toute attente, la stratégie $(I_t(\cdot))_{t \geq 1}$ de l'adversaire exécutée sans la preuve \uparrow est aveugle et même iid! Le terme $\sqrt{TK \ln K}$ provient de la différence entre $\min_{1 \leq i \leq K} \mathbb{E} \left[\sum_{t=1}^T l_{i,t} \right]$ et $\mathbb{E} \left[\min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} \right]$. Dès lors, minimiser le regret (externe) en information parfaite est aussi simple face à un adversaire antagoniste que face à une suite iid.

Rem 2: On peut aussi prouver une borne inférieure non-asymptotique de la forme: $\forall K \geq 1, \forall T \geq \square \ln K, R_{K,T}^{\text{minimum}} \geq \square \sqrt{TK \ln K}$ [preuve via le lemme de Ferns par ex].

2) Borne inférieure dans le cas bandits: \sqrt{TK}

Le problème de bandits est plus difficile qu'en information parfaite. On va en effet prouver une borne inférieure plus grande (de l'ordre de \sqrt{TK} contre $\sqrt{TK \ln K}$ tout à l'heure). On retrouve à un terme $\sqrt{\ln K}$ près la borne sup $\sqrt{TK \ln K}$ sur le (pseudo-)regret de l'ENA avec pertes estimées $\tilde{l}_{i,t}$.

Proposition: Pour tous $K \geq 2$ et $T \geq K/2$,

$$\begin{aligned} & \inf_{(P_t(\cdot))_{t \geq 1}} \sup_{(I_t(\cdot))_{t \geq 1}} \mathbb{E} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} \right] \\ & \geq \inf_{(P_t(\cdot))_{t \geq 1}} \sup_{\substack{I_1, \dots, I_T \in [0,1]^K}} \left\{ \mathbb{E} \left[\sum_{t=1}^T l_{I_t, t} \right] - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i,t} \right\} \\ & \geq \frac{1}{2\sqrt{12}} \sqrt{TK} \end{aligned}$$

où le sup est pris sur tous les adversaires t.q. $0 \leq l_{i,t} \leq 1$.

Preuve: Remarquons tout d'abord qu'on peut toujours supposer que

$\forall t \geq 1, I_t = \Psi(P_t, U_t)$, où $U_1, \dots, U_T \stackrel{iid}{\sim} \mathcal{U}([0,1])$, sans changer la loi de la suite (I_1, \dots, I_T) .

(Il suffit par ex de définir $I_t = i$ ssi $U_t \in \left[\sum_{j=1}^{i-1} P_{j,t}, \sum_{j=1}^i P_{j,t} \right]$)

Ramenons-nous comme précédemment à des pertes déterministes puis stochastiques :

$$\inf_{(p_t(\cdot))_{t \geq 1}} \sup_{(l_t(\cdot))_{t \geq 1}} \mathbb{E}_{\mathbb{P}_A} [\text{regret}(T)] \geq \inf_{(p_t(\cdot))_{t \geq 1}} \sup_{\underline{l}_1, \dots, \underline{l}_T \in [0,1]^K} \mathbb{E}_{\mathbb{P}_A} [\text{regret}(T)]$$

$$\geq \inf_{(p_t(\cdot))_{t \geq 1}} \sup_{\substack{\mathbb{Q}^* \text{ proba} \\ \text{sur } [0,1]^{KT}}} \int \mathbb{E}_{\mathbb{P}_A} [\text{regret}(T)] d\mathbb{Q}^*(\underline{l}_1, \dots, \underline{l}_T)$$

$$\stackrel{\text{Fubini}}{=} \inf_{(p_t(\cdot))_{t \geq 1}} \int \mathbb{E}_{\mathbb{Q}^*} [\text{regret}(T)] d\mathbb{P}_A(u_1, \dots, u_T)$$

La randomisation du statisticien provient de la suite $(u_1, \dots, u_T) \sim \mathcal{U}([0,1])^{\otimes T}$

$$\geq \inf_{(p_t(\cdot))_{t \geq 1}} \inf_{u_1, \dots, u_T \in [0,1]} \mathbb{E}_{\mathbb{Q}^*} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right]$$

$I_t = \Psi(p_t, u_t)$ est uniquement fonction du passé (via p_t) si (u_1, \dots, u_T) est fixé.

$$\geq \inf_{(I_t(\cdot))_{t \geq 1} \text{ prévisibles}} \mathbb{E}_{\mathbb{Q}^*} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right]$$

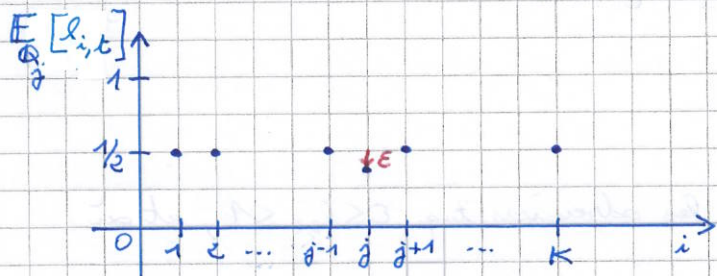
I_t fonction mesurable du passé $(l_{I_s, s})_{s \leq t-1}$

On s'est donc ramené à des algorithmes déterministes. Soit $(I_t(\cdot))_{t \geq 1}$ une telle suite de fonctions.

On va prendre $\mathbb{Q}^* = \frac{1}{K} \sum_{j=1}^K \mathbb{Q}_j$ où $\mathbb{Q}_j \in \mathcal{D}_1^+([0,1]^{KT})$ est telle que :

- sous \mathbb{Q}_j , $l_{i, t}$ indépendants, $1 \leq i \leq K$, $1 \leq t \leq T$
- $l_{i, t} \sim \begin{cases} \mathcal{B}(\frac{1}{2}) & \text{si } i \neq j \\ \mathcal{B}(\frac{1}{2} - \varepsilon) & \text{si } i = j \end{cases}$, avec $\varepsilon > 0$ à choisir ultérieurement.

(en d'autres termes : $\mathbb{Q}_j = \left(\bigotimes_{i=1}^K \mathcal{B}(\frac{1}{2} - \varepsilon \mathbb{1}_{i=j}) \right)^{\otimes T}$)



Sous \mathbb{Q}_j , le bras j est optimal et meilleur que les autres à ε près. On a :

$$\mathbb{E}_{\mathbb{Q}_j} \left[\sum_{t=1}^T l_{I_t, t} \right] = \sum_{t=1}^T \sum_{i=1}^K \mathbb{E}_{\mathbb{Q}_j} [l_{i, t} \mathbb{1}_{I_t=i}] = \sum_{t=1}^T \sum_{i=1}^K \mathbb{E}_{\mathbb{Q}_j} \left[\left(\frac{1}{2} - \varepsilon \mathbb{1}_{i=j} \right) \mathbb{1}_{I_t=i} \right]$$

conditionnement / $(l_s)_{s \leq t-1}$ car I_t prévisible. (1)

$$\text{D'où } E_{\mathbb{Q}_j} \left[\sum_{t=1}^T l_{I_t, t} \right] = \frac{T}{2} - \varepsilon \sum_{t=1}^T \mathbb{Q}_j(I_t = j)$$

Par ailleurs,

$$E_{\mathbb{Q}_j} \left[\min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] \leq \min_{1 \leq i \leq K} E_{\mathbb{Q}_j} \left[\sum_{t=1}^T l_{i, t} \right] = \frac{T}{2} - T\varepsilon$$

En conséquent,

$$\begin{aligned} E_{\mathbb{Q}_j} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] &\geq \frac{T}{2} - \varepsilon \sum_{t=1}^T \mathbb{Q}_j(I_t = j) - \left(\frac{T}{2} - T\varepsilon \right) \\ &= T\varepsilon \left(1 - \frac{1}{T} \sum_{t=1}^T \mathbb{Q}_j(I_t = j) \right) \end{aligned}$$

$$\text{D'où } E_{\mathbb{Q}^*} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] \geq T\varepsilon \left(1 - \frac{1}{T} \sum_{t=1}^T \frac{1}{K} \sum_{j=1}^K \mathbb{Q}_j(I_t = j) \right) \quad (*)$$

Reste à choisir ε suffisamment petit pour que, à chaque date t , le statisticien ne soit pas capable (quel que soit son algo) de détecter le bras optimal j avec proba suffisamment proche de 1, en moyenne sur $j \in \{1, \dots, K\}$.

↳ problème de test d'hypothèses multiples

↳ \triangle Il faut bien prendre en compte l'information de type "bandits".

Rappel: inégalité de Pinsker

Soit P et Q deux probabilités sur un même espace mesurable (Ω, \mathcal{G}) .

$$\text{Alors: } \|P - Q\|_{TV} \leq \sqrt{\frac{KL(P, Q)}{2}}$$

$$\text{où } \|P - Q\|_{TV} := \sup_{A \in \mathcal{G}} |P(A) - Q(A)| \text{ et } KL(P, Q) := \begin{cases} \int dP \log \frac{dP}{dQ} & \text{si } P \ll Q \\ \infty & \text{sinon} \end{cases}$$

Soit $t \in \{1, \dots, T\}$ fixé. D'après l'inégalité de Pinsker, on a, pour tout $j \in \{1, \dots, K\}$,

$$\mathbb{Q}_j(I_t = j) \leq \mathbb{Q}_0(I_t = j) + \sqrt{\frac{1}{2} KL(\mathbb{Q}_0^{I_t}, \mathbb{Q}_j^{I_t})}, \text{ où } \mathbb{Q}_0^{I_t} \text{ et } \mathbb{Q}_j^{I_t} \text{ désignent}$$

les lois de I_t sous \mathbb{Q}_0 et \mathbb{Q}_j respectivement,

et où $\mathbb{Q}_0 = \mathcal{B}(1/2)^{\otimes TK}$ rend les pertes $l_{i, t}$ iid $\mathcal{B}(1/2)$.

On en déduit:

$$\frac{1}{K} \sum_{j=1}^K \mathbb{Q}_j(I_t = j) \stackrel{\text{ Jensen }}{\leq} \frac{1}{K} + \sqrt{\frac{1}{2} \frac{1}{K} \sum_{j=1}^K KL(\mathbb{Q}_0^{I_t}, \mathbb{Q}_j^{I_t})}$$

Majorons $KL(Q_0^{I_t}, Q_j^{I_t})$:

$$KL(Q_0^{I_t}, Q_j^{I_t}) \leq KL(Q_0^{(I_{s,o})_{s \leq t-1}}, Q_j^{(I_{s,o})_{s \leq t-1}}) \text{ car } I_t \text{ est fonction mesurable de } (I_{s,o})_{1 \leq s \leq t-1} \leftarrow \text{bandits}$$

"chain rule for relative entropy" \rightarrow

$$= \sum_{s=1}^{t-1} \int dQ_0^{(I_{s,o})_{s \leq s-1}} KL(Q_0^{I_{s,o} | (I_{s,o})_{s \leq s-1}}, Q_j^{I_{s,o} | (I_{s,o})_{s \leq s-1}})$$

loi conditionnelle de $I_{s,o}$ sachant $(I_{s,o})_{s \leq s-1}$ sous Q_0 . idem pour Q_j

$$= \sum_{s=1}^{t-1} \int dQ_0^{(I_{s,o})_{s \leq s-1}} KL\left(\mathcal{B}\left(\frac{1}{2}\right), \mathcal{B}\left(\frac{1}{2} - \varepsilon \mathbb{1}_{I_s=j}\right)\right)$$

car I_s est fonction mesurable de $(I_{s,o})_{s \leq s-1}$

$$\leq \sum_{s=1}^{t-1} \int dQ_0^{(I_{s,o})_{s \leq s-1}} 6\varepsilon^2 \mathbb{1}_{I_s=j} \quad \boxed{\text{si } \varepsilon \leq 1/\sqrt{6}} \quad (\text{rel localement quadratique})$$

$$= 6\varepsilon^2 \sum_{s=1}^{t-1} Q_0(I_s=j)$$

$$\text{Donc } \frac{1}{K} \sum_{j=1}^K KL(Q_0^{I_t}, Q_j^{I_t}) \leq 6\varepsilon^2 \sum_{s=1}^{t-1} \frac{1}{K} \sum_{j=1}^K Q_0(I_s=j) \leq \frac{6\varepsilon^2 T}{K}$$

$$\leq \frac{1}{2} \text{ pour } \varepsilon = \sqrt{\frac{K}{12T}} \quad (\text{N.B. } \varepsilon \leq \frac{1}{\sqrt{6}} \text{ si } T \geq \frac{K}{2})$$

En injectant cette majoration pour tout $t \in \{1, \dots, T\}$ dans (*), il vient :

$$\mathbb{E}_{Q^*} \left[\sum_{t=1}^T l_{I_t, t} - \min_{1 \leq i \leq K} \sum_{t=1}^T l_{i, t} \right] \geq TE \left(1 - \frac{1}{2}\right) \text{ dès que } T \geq \frac{K}{2}$$

$$= \frac{\sqrt{TK}}{2\sqrt{12}}$$