

L'algorithme UCB, et comment aller plus loin

I. Rappel: Borne de Chernoff-Hoeffding pour les variables bornées

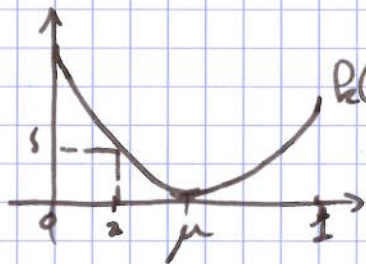
Soit X_1, \dots, X_n i.i.d $v \in \mathcal{U}_s(0,1)$ - on note $\mu = \mathbb{E}_v[X_1] = \mathbb{E}(v)$

Log-Laplace: $\phi_v(\lambda) = \mathbb{E}_v[e^{\lambda X_1}] \leq \exp(\lambda \mu + \lambda^2 \sigma^2)$

\rightarrow si $x < \mu$, on utilise Markov + optimise en λ et on obtient

$$\mathbb{P}_v(\bar{X}_n < x) \leq \exp(-n \text{KL}(x; \mu)) \quad (1)$$

où $\text{KL} : (p, q) \rightarrow p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} = \text{divergence sur } [0,1]$



\leftarrow point de vue "probable": \bar{X}_n est probable et doit s'éloigner de μ par $\text{KL}(\cdot, \mu)$
posons $\delta = \text{KL}(x; \mu)$.

(1) peut se ré-écrire $\mathbb{P}_v(\bar{X}_n < \mu - \delta, \text{KL}(\bar{X}_n, \mu) > \delta) \leq \exp(-n\delta)$

et donc

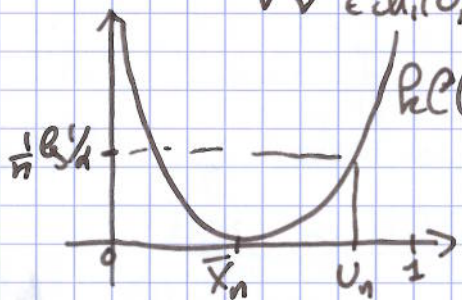
$$\forall \varepsilon > 0, \mathbb{P}_v(\bar{X}_n < \mu - \varepsilon, n \text{KL}(\bar{X}_n, \mu) > \varepsilon) \leq \exp(-\varepsilon)$$

D'où la borne supérieure de confiance

$$U_n = \max \{ \mu' \in (0,1) : n \text{KL}(\bar{X}_n, \mu') \leq \log \frac{1}{\alpha} \}$$

Elle vérifie en effet

$$\forall v \in \mathcal{U}_s(0,1), \mathbb{P}_v(U_n \geq \mathbb{E}(v)) \geq 1 - \alpha$$



\leftarrow point de vue "stats":
la moyenne est dans un voisinage de \bar{X}_n par la 2-énergie de KL

Relaxation de Hoeffding:

d'égalité de Pinsker $\text{KL}(p, q) \geq 2(p-q)^2$ (bon si $p \pm q \leq 1/2$)

$$\rightarrow \forall x < \mu, \mathbb{P}_v(\bar{X}_n < x) \leq \exp(-2n(x-\mu)^2)$$

$$\mathbb{P}_v(\bar{X}_n < \mu, 2n(\bar{X}_n - \mu)^2 > \varepsilon) \leq \exp(-\varepsilon)$$

$$\tilde{U}_n = \max \{ \mu' \in (0,1) : 2n(\bar{X}_n - \mu') \leq \log \frac{1}{\alpha} \} = \bar{X}_n + \sqrt{\frac{\log \frac{1}{\alpha}}{2n}}$$

Bonne approximation si $n \gg 1$ et n grand

II Application à l'allocation dynamique de ressources: Algorithme UCB

On présente dans cette partie les résultats principaux de

[1] Auer, Coa-Bianchi, Fischer (2002)

Finite Time Analysis of the Multiarmed Bandit Problem

1. Rappel du cadre:

K bras (resp. $(X_{a,n})_{1 \leq a \leq K, n \geq 1}$) famille de v.a. indépendantes, $X_{a,n} \sim \mathcal{U}_{[0,1]}$

On note $\mu_a = E(\mathcal{U}_a) = E[X_{a,1}]$, $\bar{X}_{a,n} = \frac{1}{n} \sum_{k=1}^n X_{a,k}$

On se définit une règle d'allocation dynamique $\pi = (\pi_t)_{t \geq 1}$

$\pi_t = (\pi_t)_{t \geq 1}$ $\pi_t: (a_1, \dots, a_{t-1}) \rightarrow a_t$

Par récurrence, cela définit les suites de variables aléatoires

$$A_t = \pi_t(A_1, Y_1, \dots, A_{t-1}, Y_{t-1})$$

$$Y_t = X_{t, N_a(t)}$$

$$N_a(t) = \sum_{s \leq t} \mathbb{1}_{\{A_s = a\}} = N_a(t-1) + \mathbb{1}_{\{A_t = a\}}$$

$$S_a(t) = \sum_{s \leq t} \mathbb{1}_{\{A_s = a\}} Y_s = \sum_{n=1}^{N_a(t)} X_{a,n}$$

$$\bar{X}_a(t) = \frac{S_a(t)}{N_a(t)} = \bar{X}_{a, N_a(t)}$$

2. Paradigme optimiste et algo UCB

"Parmi tous les environnements suffisamment raisonnables, fais comme si tu te trouvais dans celui qui t'est le plus favorable"

→ Politique UCB



$$\rightarrow A_t = \arg \max_{1 \leq a \leq K} \bar{X}_a(t) + \sqrt{\frac{\ln(t)}{2N_a(t)}}$$

(Le facteur $\frac{1}{2}$ est nécessaire dans la preuve de [1]), interprétation $\alpha = \frac{1}{t}$

Théorème (1.1): sans la politique UCB,

pour tout bras a $t \mu_a < \mu^* = \max_{1 \leq k \leq K} \mu_k$ pour tout $T \geq 1$

$$E[N_a(T)] \leq \frac{8 \ln(T)}{(\mu^* - \mu_a)^2} + \frac{\pi^2}{3} + 1$$

3. Preuve du théorème

on pose $n_T = \lceil \frac{8 \ln T}{(\mu^* - \mu_a)^2} \rceil$ et $a^* \in \arg\max_{1 \leq b \leq K} \mu_b$

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq n_T + \sum_{t=K}^{T-1} \mathbb{P}(A_{t+1}=a, N_a(t) \geq n_T) \\ &\leq n_T + \sum_{t=K}^{T-1} \mathbb{P}(L_a(t) > \mu_a) \cdot \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}(U_{a^*}(t) < \mu_{a^*}) \\ &\quad + \sum_{t=K}^{T-1} \mathbb{P}(A_{t+1}=a, N_a(t) \geq n_T, L_a(t) \leq \mu_a, U_a(t) > \mu_{a^*}) \end{aligned}$$

où $[L_a(t), U_a(t)] = \bar{X}_a(t) \pm \sqrt{\frac{4 \ln t}{2N_a(t)}}$

- $\sum_{t=1}^T \mathbb{P}(L_a(t) > \mu_a) = \sum_{t=1}^T \sum_{n=1}^t \mathbb{P}(L_a(t) > \mu_a, N_a(t) = n)$
 $= \sum_{n=1}^T \sum_{t=n}^T \mathbb{P}(\bar{X}_n > \mu_a + \sqrt{\frac{4 \ln t}{2n}}, N_a(t) = n)$
 $\leq \sum_{n=1}^T \sum_{t=n}^T \frac{1}{t^4} \leq \sum_{n=1}^T t \cdot \frac{1}{t^4} \leq \sum_{n=1}^{\infty} \frac{1}{t^2} \leq \frac{\pi^2}{6}$

et, de même, $\sum_{n=1}^T \mathbb{P}(U_{a^*}(t) < \mu_{a^*}) \leq \frac{\pi^2}{6}$

• en outre, si $L_a(t) < \mu_a$, $U_{a^*}(t) > \mu_{a^*}$ et $A_{t+1}=a$ alors

$$\mu_{a^*} < U_{a^*}(t) \leq U_a(t) = L_a(t) + 2\sqrt{\frac{4 \ln t}{2N_a(t)}} < \mu_a + \sqrt{\frac{8 \ln t}{N_a(t)}}$$

et donc $N_a(t) < \frac{8 \ln t}{(\mu_{a^*} - \mu_a)^2} \leq n_T$

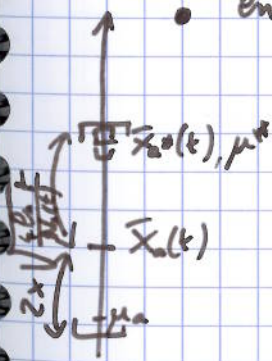
et donc $\forall t, \mathbb{P}(A_{t+1}=a, N_a(t) \geq n_T, L_a(t) \leq \mu_a, U_a(t) > \mu_{a^*}) = 0$
 d'où le résultat.

4. Commentaires :

$$\mathbb{E}[N_a(T)] \leq \frac{16 \ln T}{2(\mu_{a^*} - \mu_a)^2} (+c)$$

Version laibléris: $\mathbb{E}[N_a(T)] \geq \frac{\ln T}{\text{Rel}(\mu_a, \mu_{a^*})} (1.0(1))$

← pour la borne de l'union
 ← parce qu'a attend temps log type
 ← Spinner: $\text{Rel}(\mu_a, \mu_{a^*}) \geq \frac{\mu_{a^*} - \mu_a}{\mu_{a^*}}$



5. Améliorations

1. Utiliser U_n et pas \tilde{U}_n :

$$A_{n+1} = \arg \max_{1 \leq k \leq K} \mu'_k \in (0, 1]: N_k(t) \text{Pr}(\bar{X}_k(t), \mu'_k) \leq f(t)$$

$$\text{avec } f(t) = B_3 t (+ 3 \log B_3 t)$$

2. Éviter la borne de l'union: inégalité auto-normalisée

$$\Pr(N_k(t) \text{Pr}(\bar{X}_k(t), \mu'_k) > f(t)) \leq e^{-f(t)/B_3 t} e^{-\beta t}$$

obtenue en prouvant que

$$\forall \varepsilon > 0, \Pr\left(\bigcup_{n \leq t} \{n \text{Pr}(\bar{X}_{a,n}, \mu'_a) > \varepsilon\}\right) \leq e^{-\varepsilon/B_3 t} e^{-\beta t}$$

Preuve: reprendre celle de L.I. par martingales

3. Décomposition plus fine:

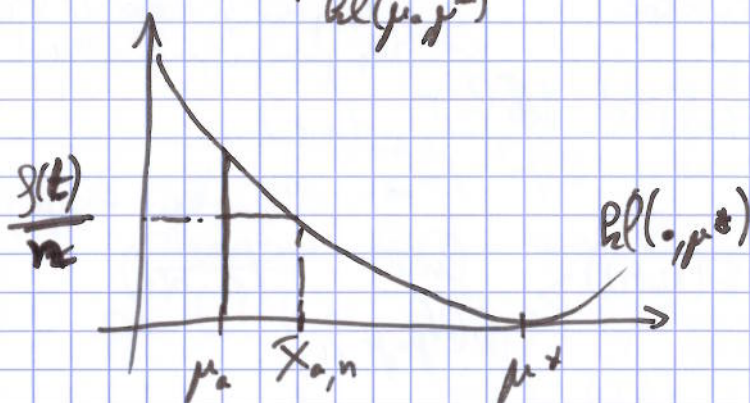
$$\{A_{n+1} = a\} \subset \{U_{a+}(t) < \mu^* \} \cup \{A_{n+1} = a, U_a(t) \geq \mu^*\}$$

$$\hookrightarrow \text{on obtient } \mathbb{E}[N_a(t)] \leq \frac{B_3 t}{\text{Pr}(\mu_a, \mu^*)} + d(\sqrt{B_3 t})$$

↑
complètement explicite

cf Gopál, Garivier, Maillard, Munos, Stoltz
Kullback-Leibler Upper Confidence Bounds
for Optimal Sequential Allocation
Annals of Statistics, 2013

Idée: d'où vient $\frac{B_3 t}{\text{Pr}(\mu_a, \mu^*)}$?



→ c'est le premier n
pour lequel $U_a(t) \geq \mu^*$
correspond à une
déviation à droite
sur $\bar{X}_{a,n}$