



Gilles Stoltz

Researcher, CNRS

Affiliate professor, HEC Paris

Paris-Saclay master in mathematics (tracks: Optimization / Data science / ALEA)

Sequential learning, sequential optimization

Final exam – Wednesday April 13, 2016

EXERCISE #1 (α, Ψ) -UCB

Consider $\left\{ \begin{array}{l} \text{a (convex) function } \Psi \\ \text{a bandit model } \mathcal{D} \end{array} \right\}$ such that all distributions $\mathcal{J} \in \mathcal{D}$ satisfy

$$\forall d \geq 0, \quad \max \left\{ \ln \mathbb{E} \left[e^{\alpha(X-EX)} \right], \ln \mathbb{E} \left[e^{\alpha(EX-X)} \right] \right\} \leq \Psi(d)$$

where $X \sim \mathcal{J}$.

For example, Hoeffding's lemma shows that for $\mathcal{D} = \mathcal{P}([0,1])$, $\Psi(d) = d^2/8$ is a suitable choice.

For all $x \geq 0$, we define $\Psi^*(x) = \sup_{d \geq 0} \{ dx - \Psi(d) \}$ and assume that Ψ^* is invertible, with inverse denoted by $(\Psi^*)^{-1}$.

We generalize UCB in the following way:

Algorithm (α, Ψ) -UCB : for a bandit problem with K arms

Parameters $\alpha > 0$ and $\Psi: [0, +\infty) \rightarrow \mathbb{R}$

For $t = 1, 2, \dots, K$: Pull each arm once, $I_t = j_t$ and get a reward $Y_t \sim \mathcal{V}_{I_t}$

For $t \geq K+1$: - Pull $I_t \in \arg \max_{j \in \{1, \dots, K\}} \left\{ \hat{\mu}_{j,t-1} + (\Psi^*)^{-1} \left(\frac{\alpha \ln t}{N_j(t-1)} \right) \right\}$

where $N_j(t-1) = \sum_{s=1}^{t-1} \mathbb{1}_{\{I_s = j\}}$ (ties broken arbitrarily)

and $\hat{\mu}_{j,t-1} = \frac{1}{N_j(t-1)} \sum_{s=1}^{t-1} Y_s \mathbb{1}_{\{I_s = j\}}$

- Get a reward Y_t drawn at random according to \mathcal{V}_{I_t} (conditionally to I_t)

We want to bound the pseudo-regret of (α, Ψ) -UCB:

$$\bar{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T Y_t\right]$$

where $\mu_j = \mathbb{E}(Y_j)$ is the expectation of $Y_j \in \mathcal{D}$
 $\mu^* = \max_{j=1, \dots, K} \mu_j$

(1) We first establish that

$$\hat{\mu}_{j,t+1} + (\Psi^*)^{-1}\left(\frac{\alpha \ln t}{N_j(t+1)}\right)$$
 is an upper confidence bound on μ_j (for all j).

(1.1) Show that for all $d \in \mathbb{R}$,

$$M_t(d) = \exp\left(d \sum_{s=1}^t (Y_s - \mu_j) \mathbb{1}_{\mathcal{I}_s = j} - \Psi(d) N_j(t)\right)$$

is a supermartingale, ie, that

$$\mathbb{E}[M_t(d) \mid \mathcal{F}_{t-1}] \leq M_{t-1}(d)$$

for a filtration (\mathcal{F}_s) to indicate.

Hint:
$$e^{d(Y_t - \mu_j) \mathbb{1}_{\mathcal{I}_t = j}} = \mathbb{1}_{\mathcal{I}_t \neq j} + \mathbb{1}_{\mathcal{I}_t = j} e^{d(Y_t - \mu_j)}$$

(1.2) Prove that $\forall \varepsilon > 0$, $\forall \ell \geq 1$, $\forall t \geq K+1$,

$$\mathbb{P}\left\{\hat{\mu}_{j,t+1} + \varepsilon \leq \mu_j \text{ and } N_j(t+1) = \ell\right\} \leq e^{-\ell \Psi^*(\varepsilon)}$$

(1.3) Deduce that
$$\mathbb{P}\left\{\hat{\mu}_{j,t+1} + (\Psi^*)^{-1}\left(\frac{\alpha \ln t}{N_j(t+1)}\right) \leq \mu_j\right\} \leq \frac{1}{t^{\alpha-1}}.$$

(2) We now establish the regret bound. Fix a suboptimal arm j .

(2.1) Explain why $I_t = j$ for $t \geq K+1$ entails one of the following, where a^* is a fixed optimal arm:

$$\hat{\mu}_{a^*, t-1} + (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_{a^*}(t-1)}\right) \leq \mu^*$$

$$\hat{\mu}_{j, t-1} - (\psi^*)^{-1}\left(\frac{\alpha \ln t}{N_j(t-1)}\right) > \mu_j$$

$$N_j(t-1) < \frac{\alpha \ln T}{\psi^*(\Delta_j/2)} \quad \text{where } \Delta_j = \mu^* - \mu_j$$

(2.2) Establish a regret bound of the form:

$$\bar{R}_T \leq \sum_{j: \Delta_j > 0} \left(\frac{\alpha \Delta_j}{\psi^*(\Delta_j/2)} \ln T + \frac{\alpha}{\alpha-2} \right).$$

(3) Can this bound be related to the one we proved for UCB on $\mathcal{D} = \mathcal{P}([0,1])$? The latter was equal to

$$\sum_{j: \Delta_j > 0} \left(\frac{8}{\Delta_j} \ln T + 2 \right).$$

EXERCISE #2.Budgeted prediction.

With fixed and known:

horizon T & loss range $[0,1]$
budget m For each round $t = 1, 2, \dots, T$:

- the decision-maker and the opponent player simultaneously pick $I_t \in \{1, \dots, N\}$ (possibly at random according to some $p_t \in \mathcal{P}(\{1, \dots, N\})$) and $\underline{l}_t = (l_{1t}, \dots, l_{Nt}) \in [0,1]^N$
- the opponent player observes I_t and p_t
- the decision-maker decides whether or not she wants to observe the \underline{l}_t ; she can do so only if she hasn't requested to perform such an observation more than $m-1$ times in total so far.

Thus the decision-maker can observe at most m vectors among the $\underline{l}_1, \dots, \underline{l}_T$; she decides which she wants to observe.

Strategy / Let's construct it step by step. We fix $\delta \in (0,1)$.

The decision-maker uses a sequence Z_1, \dots, Z_T of iid random variables distributed according to a Bernoulli distribution with parameter $\varepsilon > 0$ to decide whether she request the observation at round t ($Z_t=1$) or not ($Z_t=0$).

(1) To respect the budget constraint we want to pick ε s.t.

$$\mathbb{P}\{Z_1 + \dots + Z_T \leq m\} \geq 1 - \delta.$$

Show that $\varepsilon_t = \frac{m}{T} - \frac{1}{T} \sqrt{\frac{m}{\delta}}$ is a suitable choice when $\delta \geq \frac{1}{m}$.
 (Resort to Chebyshev's inequality:
 $\forall \varepsilon > 0, \mathbb{P}[X - \mathbb{E}X > \varepsilon] \leq \frac{\text{Var}(X)}{\varepsilon^2}$.)

(2) We define $\hat{l}_{jt} = \frac{l_{jt}}{\varepsilon} z_t$ as our estimator for l_{jt} .
 Show that for a well-chosen filtration $(\mathcal{F}_s)_{s \geq 0}$ we have:

$$\forall j, \forall t, \quad \mathbb{E}[\hat{l}_{jt} | \mathcal{F}_{t-1}] = l_{jt}.$$

(3) A lemma: shows that for all $\eta > 0$ and all non-negative numbers $u_{jt} \geq 0$, $j \in \{1, \dots, N\}$ and $t \in \{1, \dots, T\}$, we have

$$\frac{\sum_{t=1}^T \sum_{i=1}^N \frac{e^{-\eta \sum_{s=1}^{t-1} u_{is}}}{\sum_{k=1, \dots, N} e^{-\eta \sum_{s=1}^{t-1} u_{ks}}} u_{it}}{\sum_{k=1, \dots, N} \sum_{t=1}^T u_{kt}} = \min_{k=1, \dots, N} \sum_{t=1}^T u_{kt} \leq \frac{\ln N}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N q_{it} u_{it}^2$$

denoted by q_{it}

Hint: use $e^{-x} \leq 1 - x + \frac{x^2}{2}$.

(4) We consider the probability distributions

$$\tilde{p}_{jt} = \frac{\exp(-\eta \sum_{s=1}^{t-1} \hat{l}_{js})}{\sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} \hat{l}_{ks})}$$

and $I_t \sim (\tilde{p}_{1t}, \dots, \tilde{p}_{Nt})$.

Show that for this strategy the pseudo-regret is controlled as

$$\mathbb{E}\left[\sum_{t=1}^T l_{I_t, t}\right] - \min_{j=1, \dots, N} \mathbb{E}\left[\sum_{t=1}^T l_{jt}\right] \leq \frac{\ln N}{\eta} + \frac{\eta T}{2\varepsilon}.$$

- (5) Conclusion: Construct a strategy that NEVER asks more than m observations and whose pseudo-regret is controlled as

$$\mathbb{E} \left[\sum_{t=1}^T \ell_{\hat{I}_t, t} \right] - \min_{j=1 \dots N} \mathbb{E} \left[\sum_{t=1}^T \ell_{jt} \right] \leq \frac{2T}{\sqrt{m}} \sqrt{\ln N} + \frac{4T}{m}$$

(Note that with $m=T$, we recover the classical \sqrt{T} bound.)

Two much more difficult questions: [! please do exercise #3 first]

- (6) Prove that for any forecaster

$$\sup_{\ell_{jt} \in [0,1]} \left\{ \mathbb{E} \left[\sum_{t=1}^T \ell_{\hat{I}_t, t} \right] - \min_j \sum_{t=1}^T \ell_{jt} \right\} \geq c \frac{T}{\sqrt{m}}$$

for some numerical constant $c > 0$.

- (7) Prove a high-probability bound on the regret

$$\sum_{t=1}^T \ell_{\hat{I}_t, t} - \min_{j=1 \dots N} \sum_{t=1}^T \ell_{jt}$$

EXERCISE #3Approachability of a convex set.

A loss function $\ell: \{1, \dots, N\} \times \{1, \dots, M\} \rightarrow \mathbb{R}^d$ is given and known to all players.

A closed convex set $\mathcal{G} \subset \mathbb{R}^d$ is fixed.

Setting:

At each round $t=1, 2, \dots$

- the decision-maker and the opponent simultaneously pick $I_t \in \{1, \dots, N\}$ and $J_t \in \{1, \dots, M\}$ possibly at random, according to distributions denoted by $p_t \in \mathcal{P}\{1, \dots, N\}$ and $q_t \in \mathcal{P}\{1, \dots, M\}$
- the decision-maker suffers a loss $\ell(I_t, J_t)$
- both players observe I_t and J_t

Aims:

The decision-maker wants to ensure that

$$\frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \rightarrow \mathcal{G} \quad \text{a.s.}$$

that is,

$$\inf_{C \in \mathcal{G}} \left\| C - \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \right\|_2 \rightarrow 0 \quad \text{a.s.}$$

The opponent player wants to prevent this convergence.

Blackwell's condition:

$$\forall q \in \mathcal{P}\{1, \dots, M\} \quad \exists p \in \mathcal{P}\{1, \dots, N\} \mid \ell(p, q) \in \mathcal{G}$$

where we use:

$$\begin{cases} \ell(p, j) = \sum_i p_i \ell(i, j) \\ \ell(p, q) = \sum_i \sum_j p_i q_j \ell(i, j) \end{cases}$$

Short-hand notation:

(1) Show that if Blackwell's condition does not hold, then the opponent player has a strategy such that $\exists \gamma > 0$

for all strategies of the decision-maker,

$$\text{as, } \liminf_{T \rightarrow \infty} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) \right\|_2 \geq \gamma$$

Hints: Show that $\exists q_0 \in \mathcal{J}[1, \dots, M] \mid \inf_{p \in \mathcal{I}[1, \dots, N]} \inf_{c \in \mathcal{C}} \|c - \ell(p, q_0)\|_2 > 0$

Explain also why $\left\| \frac{1}{T} \sum_{t=1}^T \ell(I_t, J_t) - \frac{1}{T} \sum_{t=1}^T \ell(p_t, q_t) \right\|_2 \rightarrow 0$ a.s.

(\hookrightarrow A list of useful results is in appendix.)

(2) We now assume that Blackwell's condition holds and design a strategy for the decision-maker:

- Play $p_1 = (1/N, \dots, 1/N)$

- For $t = 2, 3, \dots$

- Compute $\bar{m}_{t-1} = \frac{1}{t-1} \sum_{s=1}^{t-1} \ell(p_s, J_s)$

- Project onto \mathcal{C} : $\bar{c}_{t-1} = \Pi_{\mathcal{C}}(\bar{m}_{t-1})$

- Pick $p_t \in \underset{p \in \mathcal{I}[1, \dots, N]}{\text{argmin}} \max_{q \in \mathcal{J}[1, \dots, M]} \langle \bar{m}_{t-1} - \bar{c}_{t-1}, \ell(p, q) \rangle$

where $\langle u, v \rangle = \sum_i u_i v_i$ denotes the inner product of \mathbb{R}^d

(2.1) Recall on a picture why

$$\forall c \in \mathcal{C} \quad \langle \bar{m}_{t-1} - \bar{c}_{t-1}, c - \bar{c}_{t-1} \rangle \leq 0$$

(2.2) Prove that $\forall q \in \mathcal{J}[1, \dots, M], \quad \langle \bar{m}_{t-1} - \bar{c}_{t-1}, \ell(p_t, q) - \bar{c}_{t-1} \rangle \leq 0$

Hint: A consequence of Son's lemma is that

$$\max_q \min_p \langle \bar{m}_{t-1} - \bar{c}_{t-1}, \ell(p, q) - \bar{c}_{t-1} \rangle = \min_p \max_q \langle \bar{m}_{t-1} - \bar{c}_{t-1}, \ell(p, q) \rangle$$

(2.3) Denote by $d_t = \inf_{c \in \mathcal{C}} \|c - \bar{m}_t\|_2$

and by $L = \max_{i,j} |\ell(i,j)|$.

Show that $d_{t+1}^2 \leq \left(1 - \frac{2}{t+1}\right) d_t^2 + \frac{4M^2}{(t+1)^2}$

(2.4) Prove that, $\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{t} \sum_{s \leq t} \ell(p_s, J_s) \right\|_2 \leq \frac{2M}{\sqrt{t}}$

(2.5) Conclude to the desired convergence.

Note: see a list of useful results in appendix.

Useful results.Hoeffding Azuma inequality:

Let $(\mathcal{F}_t)_{t \geq 0}$ be a filtration and let $(X_t)_{t \geq 1}$ be a sequence of adapted random variables such that

$$\forall t \geq 1, \quad \exists Z_{t-1} \text{ } \mathcal{F}_{t-1}\text{-measurable,}$$

$$\exists (a_t, b_t) \in \mathbb{R}^2$$

$$| \quad \mathcal{F}_{t-1} + a_t \leq X_t \leq \mathcal{F}_{t-1} + b_t \quad \text{a.s.}$$

Then $\forall \delta \in (0, 1)$,

$$\mathbb{P} \left\{ \sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}] \leq \sqrt{\frac{\sum_{t=1}^T (b_t - a_t)^2}{2} \ln \frac{1}{\delta}} \right\} \geq 1 - \delta$$

Sion's lemma

(as seen in the exercises):

Let X, Y be two convex sets

with X metric and compact,

and let $f: X \times Y \rightarrow \mathbb{R}^d$ be

- bounded,
- uniformly continuous,
- concave in the second argument,
- convex in the first argument,

$$\text{then} \quad \inf_{x \in X} \sup_{y \in Y} f(x, y) = \sup_{y \in Y} \inf_{x \in X} f(x, y).$$