# Part 1: $\sqrt{KT}$ distribution-free regret bounds for stochastic bandits

An exercise of the homework is about proving:

Distribution-free (ie, uniform) lower bounds.

We prove:    For all $K \geq 2$ and $T \geq K/5$

$$\inf_{\substack{\text{strategies} \\ \psi}} \quad \sup_{\substack{\nu_1, \ldots, \nu_K \\ \text{in } \mathcal{P}([0,1])}} \quad R_T \geq \frac{1}{20} \sqrt{TK}.$$

and even :
$\sup$ over $\nu_1 \ldots \nu_K$
being Bernoulli distributions

We saw that UCB enjoyed a distribution-free regret bound of order $\sqrt{TK \ln T}$, but the $\sqrt{\ln T}$ is unnecessary. The optimal (minmax) distribution-free regret bound for bounded stochastic bandits is of order $\sqrt{TK}$.

We now discuss an algorithm achieving this optimal order of magnitude; it is called MOSS and is a variation on UCB, with a smaller / more careful exploration bonus.

The MOSS strategy    (Minimax Optimal Strategy in the Stochastic case of bandit problems)

Index policy relying on

$$U_a(t) = \hat{\mu}_a(t) + \sqrt{\frac{1}{2N_a(t)} \ln_+\left(\frac{t}{K N_a(t)}\right)}$$

for $t \geq K$,

and where $\ln_+ = \max\{\ln, 0\}$

That is:    For $t = 1, \dots K$:       pull arm $A_t = t$

For $t \geq K+1$:    pull arm $A_t \in \underset{a=1\dots K}{\arg\max} \; U_a(t-1)$

Difference to UCB:    we replace the exploration bonus

$$\sqrt{2\ln t / N_a(t)} \qquad \text{by} \qquad \sqrt{\ln_+\left(\frac{t}{K N_a(t)}\right) / (2 N_a(t))}$$

↳ no exploration after a was pulled sufficiently often ($t/K$ times)

We prove a distribution-free bound:

Theorem:    MOSS is such that $\underset{\substack{\nu_1 \dots \nu_K \\ \text{distributions over } [0,1]}}{\sup} \overline{R}_T \leq K-1 + 45\sqrt{KT}$

(the constant 45 can be improved)
↓
but indeed minimax optimal as its name indicates!

Open question:    Take inspiration from the MOSS proof to write a better (more direct) proof for UCB!

**Proof.**  __First step:__  $U_{a^*}(t-1) \leq U_{A_t}(t-1)$  by definition of $A_t$ as an argmax

thus  $R_T = \sum_{t=1}^{T} \mathbb{E}[\mu^* - \mu_{A_t}]$

$\qquad \leq K-1 + \underbrace{\sum_{t=K+1}^{T} \mathbb{E}[\mu^* - U_{a^*}(t-1)]}_{\substack{\text{if we played} \\ \text{each arm once in the} \\ \text{first K steps, and} \\ \text{at most K-1 were suboptimal}}} + \underbrace{\sum_{t=K+1}^{T} \mathbb{E}[U_{A_t}(t-1) - \mu_{A_t}]}_{}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \leq \sqrt{KT} + \sum_{t=K+1}^{T} \mathbb{E}\left[\left(U_{A_t}(t-1) - \mu_{A_t} - \sqrt{K/T}\right)^+\right]$

__Second step:__
__Control of__ each  $\mathbb{E}[\mu^* - U_{a^*}(t)]$  term  by  $20\sqrt{K/t}$

We write  $\mathbb{E}[\mu^* - U_{a^*}(t)]$
for $t \geq K$

$\qquad \leq \mathbb{E}\left[(\mu^* - U_{a^*}(t))^+\right]$

$\qquad \leq \sum_{\ell=0}^{+\infty} \mathbb{E}\left[(\mu^* - U_{a^*}(t))^+ \mathbb{1}_{\{N_{a^*}(t) \in [x_{\ell+1}, x_\ell]\}}\right]$   where $x_\ell = \beta^{-\ell} t/K$ for some fixed $\beta > 1$ and $\ell = 0, 1, 2, \ldots$

$\qquad\qquad + \mathbb{E}\left[(\mu^* - U_{a^*}(t))^+ \mathbb{1}_{\{N_1(t) \geq t/K\}}\right]$

Now,  $U_{a^*}(t) = \hat{\mu}_{a^*}(t) + \begin{cases} 0 & \text{if } N_1(t) \geq t/K \\[4mm] \sqrt{\dfrac{1}{2N_{a^*}(t)} \ln\left(\dfrac{t}{K N_{a^*}(t)}\right)} \geq \sqrt{\dfrac{1}{2x_\ell} \ln\left(\dfrac{t}{K x_\ell}\right)} \\ \qquad\qquad \text{if } N_{a^*}(t) \in [x_{\ell+1}, x_\ell] \end{cases}$   $\left.\begin{array}{c} \\ \\ \\ \\ \end{array}\right\}$ denoted by $\varepsilon_\ell$

Therefore,  $\mathbb{E}[\mu^* - U_{a^*}(t)]$
(*)

$\qquad \leq \mathbb{E}\left[(\mu^* - \hat{\mu}_{a^*}(t))^+ \mathbb{1}_{\{N_{a^*}(t) \geq t/K\}}\right] + \sum_{\ell=0}^{+\infty} \mathbb{E}\left[(\mu^* - \hat{\mu}_{a^*}(t) - \varepsilon_\ell)^+ \mathbb{1}_{\{N_{a^*}(t) \in [x_{\ell+1}, x_\ell]\}}\right]$

__Lemma:__  $\mathbb{E}\left[(\mu^* - \hat{\mu}_{a^*}(t) - \varepsilon)^+ \mathbb{1}_{\{N_{a^*}(t) \geq n_0\}}\right] \leq \dfrac{1}{\sqrt{n_0}} e^{-2n_0 \varepsilon^2}$

__Proof of the lemma:__
$\qquad \mathbb{E}\left[(\mu^* - \hat{\mu}_{a^*}(t) - \varepsilon)^+ \mathbb{1}_{\{N_{a^*}(t) \geq n_0\}}\right]$

$\qquad = \int_0^{+\infty} \mathbb{P}\left\{\mu^* - \hat{\mu}_{a^*}(t) - \varepsilon \geq u \quad \& \quad N_{a^*}(t) \geq n_0\right\} du$

$$= \int_0^{+\infty} \mathbb{P}\left\{ Z_t^* \geq (\varepsilon + u) N_a^*(t) \quad \& \quad N_a^*(t) \geq n_0 \right\} du$$

where

$$Z_t^* = N_a^*(t) \left( \mu^* - \hat{\mu}_{a^*}(t) \right) = \sum_{s=1}^{t} \left( \mu^* - Y_s \right) \mathbb{1}_{\{A_t = a^*\}}$$

is a martingale,

and for all $x \in \mathbb{R}$,

$$S_{x,t} = e^{x Z_t^* - x^2/8 \, N_a^*(t)}$$

is a supermartingale.

(margin, right side) we did something similar for UCB

Thus, by Markov-Chernoff, we continue the bounding as, for $x > 0$:

$$= \int_0^{+\infty} \mathbb{P}\left\{ e^{x Z_t^* - x^2/8 \, N_a^*(t)} \geq \exp\left( N_a^*(t) \left( x(\varepsilon + u) - \frac{x^2}{8} \right) \right) \quad \& \quad N_a^*(t) \geq n_0 \right\} du$$

$$\leq \int_0^{+\infty} \sum_{\ell = n_0}^{+\infty} e^{-2\ell(\varepsilon + u)^2} \, \mathbb{E}\left[ S_{4(\varepsilon + u), t} \, \mathbb{1}_{\{N_a^*(t) = \ell\}} \right] du$$

we pick,
$x = 4(\varepsilon + u)$
so that $x(\varepsilon + u) - \frac{x^2}{8} = 2(\varepsilon + u)^2$

$\not\!\!\forall$ independent of $\ell$, which will be useful in other proofs!

we
we
$e^{-2\ell(\varepsilon + u)^2}$
$\leq e^{-2n_0(\varepsilon + u)^2}$
$\leq e^{-2n_0 \varepsilon^2} e^{-2n_0 u^2}$

$$\leq e^{-2n_0 \varepsilon^2} \int_0^{+\infty} e^{-2n_0 u^2} \, \mathbb{E}\left[ S_{4(\varepsilon + u), t} \, \mathbb{1}_{\{N_a^*(t) \geq n_0\}} \right] du$$

where $\mathbb{E}\left[ S_{4(\varepsilon + u), t} \right] \leq 1$

All in all, $\mathbb{E}\left[ \left( \mu^* - \hat{\mu}_{a^*}(t) - \varepsilon \right)^+ \mathbb{1}_{\{N_a^*(t) \geq n_0\}} \right]$

$$\leq e^{-2n_0 \varepsilon^2} \int_0^{+\infty} e^{-2n_0 u^2} \, du = e^{-2n_0 \varepsilon^2} \sqrt{\frac{\pi}{8 n_0}}$$

$$\leq e^{-2n_0 \varepsilon^2} / \sqrt{n_0}$$

where $\sqrt{\frac{\pi}{8}} \leq 1$

integral of a Gaussian density, up to the normalization factor

$\hookrightarrow \sigma^2 = \frac{n_0}{4}$ in $\int_0^{+\infty} \frac{1}{\sigma \sqrt{2\pi}} e^{-u^2/2\sigma^2} \, du = \frac{1}{2}$

Substituting the lemma in (*):

$$\mathbb{E}\big[\mu^* - U_{a^*}(t)\big] \leq \sqrt{\tfrac{K}{t}} + \sum_{\ell=0}^{+\infty} \underbrace{\frac{1}{\sqrt{x_{\ell+1}}} e^{-2x_{\ell+1}\varepsilon_\ell^2}}$$

$$\underbrace{\frac{1}{\sqrt{x_{\ell+1}}} \exp\Big(-2x_{\ell+1} \underbrace{\frac{1}{2x_{\ell+1}}}_{1/\beta} \underbrace{\ln\big(\tfrac{t}{Kx_\ell}\big)}_{\ln\beta^\ell \,=\, \ell\ln\beta}\Big)}$$

where

$$\sum_{\ell \geq 0} \beta^{\ell(1/2 - 1/\beta)} \text{ is } < +\infty$$

$$\qquad = \sqrt{\tfrac{K}{t}}\,\beta^{(\ell+1)/2}\exp\Big(-\tfrac{\ell}{\beta}\ln\beta\Big)$$

as soon as $\beta \in (1,2)$

$$= \sqrt{\tfrac{K}{t}}\,\beta^{1/2 + \ell(1/2 - 1/\beta)}$$

E.g., for $\beta = \tfrac{3}{2}$,

$$\sum_{\ell=0}^{+\infty}\Big(\tfrac{3}{2}\Big)^{1/2 + \ell(1/2 - 2/3)} = \sqrt{\tfrac{3}{2}}\sum_{\ell=0}^{+\infty}\alpha^{+\ell} = \frac{1}{1-\alpha}\sqrt{\tfrac{3}{2}} \leq 19$$

where $\alpha = \Big(\tfrac{3}{2}\Big)^{1/2 - 2/3} \in (0,1)$

All in all:   we obtain a   $\sqrt{\tfrac{K}{t}} + 19\sqrt{\tfrac{K}{t}} = 20\sqrt{\tfrac{K}{t}}$ bound, as claimed.

Third step:

$$\sum_{t=K+1}^{T} \mathbb{E}\Big[\big(U_{A_t}(t-1) - \mu_{A_t} - \sqrt{\tfrac{K}{T}}\big)^+\Big] \text{ is } \leq 4\sqrt{KT}$$

$$= \sum_{t=K}^{T-1} \mathbb{E}\Big[\big(U_{A_{t+1}}(t) - \mu_{A_{t+1}} - \sqrt{\tfrac{K}{T}}\big)^+\Big]$$

We decompose the expectations of interest according to the $\{A_{t+1} = a\}$ and $\{N_a(t) = \ell\}$:

$$\sum_{t=K}^{T-1} \mathbb{E}\Big[\big(U_{A_{t+1}}(t) - \mu_{A_{t+1}} - \sqrt{\tfrac{K}{T}}\big)^+\Big] = \sum_{a=1}^{K}\sum_{\ell=1}^{T}\sum_{t=K}^{T-1} \mathbb{E}\Big[\big(U_a(t) - \mu_a - \sqrt{\tfrac{K}{T}}\big)^+ \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}}\Big]$$

We now use

$$\big(U_a(t) - \mu_a - \sqrt{\tfrac{K}{T}}\big)^+ \leq \big(\hat\mu_a(t) - \mu_a - \sqrt{\tfrac{K}{T}}\big)^+ + \begin{cases} 0 & \text{if } N_a(t) \geq \tfrac{t}{K} \\ \sqrt{\tfrac{1}{2N_a(t)}\ln\big(\tfrac{t}{KN_a(t)}\big)} & \end{cases}$$

$$\text{also smaller than } \sqrt{\tfrac{1}{2N_a(t)}\ln\big(\tfrac{T}{KN_a(t)}\big)} \text{ if } N_a(t) < \tfrac{t}{K}$$

and get therefore the upper bound

$$\sum_{a=1}^{K}\sum_{\ell=1}^{T}\sum_{t=K}^{T-1} \mathbb{E}\Big[\big(\hat\mu_a(t) - \mu_a - \sqrt{\tfrac{K}{T}}\big)^+ \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}}\Big]$$

$$+ \sum_{a=1}^{K}\sum_{\ell=1}^{\lfloor T/K \rfloor} \sqrt{\tfrac{1}{2\ell}\ln\big(\tfrac{T}{K\ell}\big)}\; \mathbb{E}\Big[\sum_{t=K}^{T-1} \mathbb{1}_{\{N_a(t)=\ell\}} \mathbb{1}_{\{A_{t+1}=a\}}\Big]$$

We will repeatedly use that

$$\forall a, \ell, \qquad \sum_{t=k}^{T-1} \mathbb{1}_{\{N_a(t) = \ell\}} \mathbb{1}_{\{A_{t+1} = a\}} \leq 1 \qquad \text{(i.e., disjoint union)}$$

if $N_a(t)$ increases by 1 whenever $a$ is played

Also,

$$\sum_{\ell=1}^{\lfloor T/k \rfloor} \sqrt{\frac{1}{2\ell} \ln\left(\frac{T}{k\ell}\right)} \leq \int_0^{\lfloor T/k \rfloor} \sqrt{\frac{1}{2x} \ln\left(\frac{T}{kx}\right)} \, dx$$

change of variable
$u = T/(kx)$

$$\leq \sqrt{\frac{T}{2k}} \int_1^{+\infty} u^{-3/2} \sqrt{\ln u} \, du$$

by the change of variable $u = e^{v^2}$

$$= \sqrt{\frac{T}{2k}} \int_0^{+\infty} 2 v^2 \, e^{-v^2/2} \, dv$$

$$= \sqrt{\pi} \sqrt{\frac{T}{k}}$$

Summarizing what we proved so far:

will show that $\leq \sqrt{\tfrac{\pi}{2}} \sqrt{T/k}$ for each $a$

$$\sum_{t=k}^{T-1} \mathbb{E}\left[\left(U_{A_{t+1}}(t) - \mu_{A_{t+1}} - \sqrt{k/t}\right)^+\right] \leq \sqrt{\pi} \sqrt{kT} + \sum_{a=1}^{K} \sum_{\ell=1}^{T} \sum_{t=k}^{T-1} \mathbb{E}\left[\left(\hat{\mu}_a(t) - \mu_a - \sqrt{k/t}\right)^+\right] \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}}$$

We resort again to $Z_{a,t} = N_a(t)\left(\hat{\mu}_a(t) - \mu_a\right)$ — martingale

and $S^{(a)}_{x,t} = e^{x Z_{a,t} - \frac{x^2}{8} N_a(t)}$ — super martingale

where $x = 4\left(\sqrt{k/t} + u\right)$

For each $a$,

$$\sum_{\ell=1}^{T} \sum_{t=k}^{T-1} \mathbb{E}\left[\left(\hat{\mu}_a(t) - \mu_a - \sqrt{k/t}\right)^+ \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}}\right]$$

$$= \int_0^{+\infty} \left(\hat{\mu}_a(t) - \mu_a - \sqrt{k/t}\right)^+$$

$$= \sum_{\ell=1}^{T} \sum_{t=k}^{T-1} \int_0^{+\infty} \mathbb{P}\left\{ x Z_{a,t} - \frac{x^2}{8} N_a(t) \geq N_a(t)\left(x(u + \sqrt{k/t}) - \frac{x^2}{8}\right) \ \& \ A_{t+1}=a \ \& \ N_a(t)=\ell \right\} \, du$$

$$\leq \sum_{\ell=1}^{T} \sum_{t=k}^{T-1} \int_0^{+\infty} \underbrace{e^{-\ell\ell(u+\sqrt{k/t})^2}}_{\leq e^{-2\ell u^2} e^{-2\ell k/t}} \ \mathbb{E}\left[ S^{(a)}_{x,t} \ \mathbb{1}_{\{A_{t+1}=a\}} \ \underbrace{\mathbb{1}_{\{N_a(t)=\ell\}}}_{\text{the sum over } t \text{ of these will be} \leq 1} \right] \, du$$

by Markov-Chernoff

issue: this depends on $t$... but can be replaced in some sense, by $S^{(a)}_{x,0} = 1$

But:  remember Doob's maximal inequality for non-negative supermartingales:

$\llcorner$ see also an alternative treatment on the next page

$$\mathbb{P}\left\{\sup_{t\geq 0} S^{(a)}_{x,t} \geq C\right\} \leq \frac{\mathbb{E}\left[S^{(a)}_{x,0}\right]}{C} = \frac{1}{C}$$

Then,

$$\sum_{\ell=1}^{T} \sum_{t=K}^{T-1} \mathbb{E}\left[\left(\hat{\mu}_a(t)-\mu_a - \sqrt{K/T}\right)^+ \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}}\right]$$

as before!  $\llcorner$

$$= \sum_{\ell=1}^{+\infty} \sum_{t=K}^{T-1} \int_0^{+\infty} \mathbb{P}\left\{S^{(a)}_{x,t} \geq e^{2\ell(u+\sqrt{K/T})^2} \text{ and } A_{t+1}=a \text{ and } N_a(t)=\ell\right\} du$$

$$\leq \sum_{\ell=1}^{T} \int_0^{+\infty} \sum_{t=K}^{T-1} \mathbb{P}\left\{\left(\sup_{6\geq 0} S^{(a)}_{x,6}\right) \geq e^{2\ell(u+\sqrt{K/T})^2} \text{ and } A_{t+1}=a \text{ and } N_a(t)=\ell\right\} du$$

$\leq$ cf. disjoint union!

$$\sum_{\ell=1}^{T} \int_0^{+\infty} \mathbb{P}\left\{\left(\sup_{6\geq 0} S^{(a)}_{x,6}\right) \geq e^{2\ell(u+\sqrt{K/T})^2}\right\} du$$

Doob's maximal inequality $\leq$

$$\sum_{\ell=1}^{T} \int_0^{+\infty} e^{-2\ell(u+\sqrt{K/T})^2} du \qquad \leq \sum_{\ell=1}^{+\infty} \frac{1}{\sqrt{\ell}} e^{-2\ell K/T}$$

$$\underbrace{\leq e^{-2\ell u^2}} \times e^{-2\ell K/T}$$

and same treatment as in the lemma of the first part of the proof

This step is concluded by calculations:

$$\sum_{\ell=1}^{T} \frac{1}{\sqrt{\ell}} e^{-2\ell K/T} \leq \int_0^T \frac{1}{\sqrt{x}} e^{-2xK/T} dx$$

$$= \sqrt{\frac{T}{2K}} \int_0^{+\infty} \frac{e^{-u}}{\sqrt{u}} du = \sqrt{\frac{T}{2K}} \int_0^{+\infty} e^{-v^2} dv = \sqrt{\frac{\pi}{2}} \sqrt{\frac{T}{K}}$$

Final bound is:  $\sqrt{\pi}\ \sqrt{KT} + K\sqrt{\frac{\pi}{2}}\sqrt{T/K} \leq 4\sqrt{KT}$

General conclusion:  Final bound given by

$$K-1 + \left(\sum_{t=K+1}^{T} 20\sqrt{\frac{K}{t}}\right) + \sqrt{KT} + 4\sqrt{KT} \leq K-1 + 5\sqrt{KT} + 20\int_0^T \sqrt{K/t}\ dt$$

$$= K-1 + 45\sqrt{KT}$$

Alternative treatment (credits to Enzo Miller) of the end of Step #3:

We were stuck at

$$\sum_{\ell=1}^{T} \sum_{t=K}^{T-1} \int_0^{+\infty} e^{-2\ell u^2} e^{-2\ell K/T} \mathbb{E}\left[ S_{xt}^{(a)} \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=\ell\}} \right] du$$

$$= \sum_{\ell=1}^{T} \int_0^{+\infty} e^{-2\ell u^2} e^{-2\ell K/T} \mathbb{E}\left[ \underbrace{\sum_{t=K}^{T-1} S_{xt}^{(a)} \mathbb{1}_{\{A_{t+1}=a\}}}_{= S_{x,T_\ell}^{(a)}} \mathbb{1}_{\{N_a(t)=\ell\}} \right] du$$

where $T_\ell$ is given by:

$$T_\ell = \inf\{t \in \{1,...,T\} : A_{t+1}=a \text{ and } N_a(t)=\ell\}$$

We should get $\mathbb{E}[S_{x,T_\ell}^{(a)}] \leq \mathbb{E}[S_{x,0}^{(a)}]=1$ from the optional stopping theorem (« théorème d'arrêt de Doob ») provided some verifications. ($T_\ell$ should be a bounded stopping time.)

# Part 2: Adversarial bandits

## Adversarial bandits.          ( Rather stated in terms of losses, than rewards!)

<u>Setting:</u>          At each round $t = 1, 2, \ldots$

    1.   The opponent and the decision-maker simultaneously

       choose $\ell_t = (\ell_{j,t})_{j \in \{1 \ldots N\}}$ and $I_t \sim p_t$, where

$$p_t \in \mathcal{P}(\{1, \ldots N\})$$

    2.   The opponent gets to see $p_t$ and $I_t$;

       the decision-maker only observes $\ell_{I_t, t}$ (her own loss).

<u>Regret:</u>          $R_T = \sum\limits_{t=1}^{T} \ell_{I_t, t} \; - \; \min\limits_{j=1 \ldots N} \sum\limits_{t=1}^{T} \ell_{j,t}$

vs.   Pseudo - regret:          $\overline{R}_T = \mathbb{E}\left[ \sum\limits_{t=1}^{T} \ell_{I_t, t} \right] - \min\limits_{j=1 \ldots N} \mathbb{E}\left[ \sum\limits_{t=1}^{T} \ell_{j,t} \right]$

↑
same definition as for stochastic
bandits, up to the conversion
of losses $\ell_{j,t}$ into rewards $M - \ell_{j,t}$
( for a well chosen bound $M$ )

Why $\mathbb{E}$?
cf. $\ell_t$ are
random variables, as
they depend on the past,
and in particular on $I_1,$
$\ldots, I_{t-1}$

We have          $\overline{R}_T \leq \mathbb{E}[R_T]$.

We actually rather shoot for high-probability bounds on $R_T$, but studying $\overline{R}_T$ will be a good warm-up!

    ⚠   In these lecture notes, I'll take $N = K$ as the

    number of components

       ↳ we used $N$ for individual sequences

       ↳ $K$ stochastic bandits

   and I alternatively took $N$ and $K$ in the next pgs...

                ( My bad... )

Adversarial bandits:          bound on $\overline{R}_T$ via exponential weights.

Key: Estimators of the losses (the unseen and the seen ones):

$$\hat{\ell}_{jt} = \frac{\ell_{I_t,t}}{p_{jt}} \mathbb{1}_{\{I_t = j\}} \qquad \text{if } p_{jt} \neq 0 \text{ (which we will assume)}$$

auxiliary randomizations of opponent + decision $j$-maker

They are (conditionally) unbiased: denoting by $\mathcal{F}_{t-1} = \sigma\begin{pmatrix} U_1 \ldots U_{t-1}, \\ \ell_1, \ldots \ell_{t-1}, \\ p_1, \ldots p_{t-1}, \\ I_1, \ldots I_{t-1} \end{pmatrix}$

the total information available at the beginning of round $t$

(of course, the decision-maker does not have that much information!),

we have:

- $\ell_t$ and $p_t$ are $\mathcal{F}_{t-1}$-measurable; the only randomness comes from the random draw of $I_t$ according to $p_t$ **thanks to $U_t$**

- $\hat{\ell}_{jt}$ can be rewritten $\hat{\ell}_{jt} = \frac{\ell_{jt}}{p_{jt}} \mathbb{1}_{\{I_t = j\}}$

so that

$$\mathbb{E}\big[ \hat{\ell}_{jt} \mid \mathcal{F}_{t-1} \big] = \frac{\ell_{jt}}{p_{jt}} \underbrace{\mathbb{E}\big[ \mathbb{1}_{\{I_t = j\}} \mid \mathcal{F}_{t-1} \big]}_{= p_{jt}} = \frac{\ell_{jt}}{p_{jt}} p_{jt} = \ell_{jt}$$

since we assumed $p_{jt} \neq 0$.

Algorithm: $p_1 = (1/N, \ldots 1/N)$ and for $t \geqslant 2$, $p_t = (p_{jt})_{j=1,\ldots N}$ is defined as

for a non-increasing sequence $(\eta_t)_{t \geqslant 2}$

$$p_{jt} = \exp\left(-\eta_t \sum_{s=1}^{t-1} \hat{\ell}_{js}\right) \Big/ \sum_{k=1}^{N} \exp\left(-\eta_t \sum_{s=1}^{t-1} \hat{\ell}_{ks}\right)$$

↳ ensures indeed that $p_{jt} \neq 0$.

the range $[0, M]$ is assumed to be known...

Theorem: The strategy above, tuned with $\eta_t = \frac{1}{M} \sqrt{\frac{\ln N}{N t}}$, is such that:

for all opponents picking losses $\ell_{jt} \in [0, M]$,

$$\overline{R}_T = \mathbb{E}\left[ \sum_{t=1}^{T} \ell_{I_t, t} \right] - \min_{i=1\ldots N} \mathbb{E}\left[ \sum_{t=1}^{T} \ell_{it} \right] \leqslant 2M \sqrt{T N \ln N}$$

The proof is based on the following lemma.

**UNBOUNDED and $\geq 0$**

**Lemma :** The exponentially weighted average strategy on losses

$$\tilde{\ell}_{jt} \in [0, +\infty[, \quad \text{ie.,}$$

**... which is why we develop a new upper bound**

with $\eta_t \downarrow$,

$$\tilde{p}_{jt} = \exp\left(-\eta_t \sum_{s=1}^{t-1} \tilde{\ell}_{js}\right) \bigg/ \sum_{k=1}^{N} \exp\left(-\eta_t \sum_{s=1}^{t-1} \tilde{\ell}_{ks}\right),$$

is such that

$$\sum_{t=1}^{T} \sum_{j=1}^{N} \tilde{p}_{jt} \tilde{\ell}_{jt} \; - \; \min_{i=1\cdots N} \sum_{t=1}^{T} \tilde{\ell}_{it} \; \leq \; \frac{\ln N}{\eta_T} + \sum_{t=1}^{T} \frac{\eta_t}{2} \sum_j \tilde{p}_{jt} \tilde{\ell}_{jt}^2.$$

**Proof :** We saw earlier in this series of lectures that the EWA strategy ( with $\eta_t \downarrow$ ) is such that

$$\forall \tilde{\ell}_{jt} \in \mathbb{R}, \quad \sum_{t,j} \tilde{p}_{jt} \tilde{\ell}_{jt} \; - \; \min_{i=1\cdots N} \sum_{t=1}^{T} \tilde{\ell}_{it} \; \leq \; \frac{\ln N}{\eta_T} + \sum_{t=1}^{T} \tilde{\delta}_t$$

where $\quad \tilde{\delta}_t = \sum_{j=1}^{N} \tilde{p}_{jt} \tilde{\ell}_{jt} + \frac{1}{\eta_t} \ln \sum_{j=1}^{N} \tilde{p}_{jt} e^{-\eta_t \tilde{\ell}_{jt}}$

We use here $\quad e^{-x} \leq 1 - x + \frac{x^2}{2} \quad \forall x \geq 0$

So that $\quad \ln \sum_j \tilde{p}_{jt} e^{-\eta_t \tilde{\ell}_{jt}} \; \leq \; \ln\left(1 - \eta_t \sum_j \tilde{p}_{jt} \tilde{\ell}_{jt} + \frac{\eta_t^2}{2} \sum_j \tilde{p}_{jt} \tilde{\ell}_{jt}^2\right)$

$$\underset{\substack{\ln(1+u) \leq u \\ \forall u \geq -1}}{\leq} \; -\eta_t \sum_j \tilde{p}_{jt} \tilde{\ell}_{jt} + \frac{\eta_t^2}{2} \sum_j \tilde{p}_{jt} \tilde{\ell}_{jt}^2$$

hence the stated bound.

**Proof (of the theorem) :** We have no control on how large the $\hat{\ell}_{jt}$ can be, and they could be very large ! So, we would not be ready to apply any bound with a remainder $M_T \ln N$ term, where $M_T$ is such that $\hat{\ell}_{jt} \in [0, M_T] \; \forall j,t \cdots$ as this $M_T$ could be even super-linear. That's why we go back to the beginning of the

proof for the fully adaptive algorithm .... The $\eta_t$ can be picked as

$\ln N / \sum\limits_{s=1}^{t-1} \delta_s$ — or in terms of the upper bounds on the $\delta_s$ ( we choose the latter version for the sake of concreteness).

$\hookrightarrow$ see below!

The lemma yields for the $\hat{\ell}_{jt}$ :

$$\sum_{t,j} p_{jt}\, \hat{\ell}_{jt} \;-\; \min_{i=1,\ldots N} \sum_{t=1}^{T} \hat{\ell}_{it} \;\le\; \frac{\ln N}{\eta_T} \;+\; \sum_{t=1}^{T} \frac{\eta_t}{2} \sum_{j} p_{jt}\, \hat{\ell}_{jt}^{\,2}$$

by definition of the $\hat{\ell}_{jt}$, this

$$= \sum_{j} \ell_{I_t,t}\, \frac{p_{jt}}{p_{jt}}\, \mathbb{1}_{\{I_t = j\}} = \ell_{I_t,t}$$

similar treatment:

$$= \sum_{j} \ell_{I_t,t}^{2}\, \frac{1}{p_{jt}}\, \mathbb{1}_{\{I_t=j\}}$$
$$\le M^2 \sum_{j} \mathbb{1}_{\{I_t=j\}} / p_{jt}$$

To simplify even further the choice of the $\eta_t$, we first take $\mathbb{E}$ of both sides:

$$\mathbb{E}\!\left[\sum_{t} \ell_{I_t,t}\right] \;-\; \mathbb{E}\!\left[\min \sum_{t} \hat{\ell}_{it}\right] \;\le\; \frac{\ln N}{\eta_T} \;+\; \frac{M^2}{2} \sum_{t=1}^{T} \eta_t \sum_{j=1}^{N} \mathbb{E}\!\left[\frac{\mathbb{1}_{\{I_t=j\}}}{p_{jt}}\right]$$

$\le \min \sum\limits_{t=1}^{T} \mathbb{E}[\hat{\ell}_{it}]$

$= 1$

$= \mathbb{E}[\ell_{it}]$ by the tower rule and the fact that $\hat{\ell}_{it}$ is conditionally unbiased.

Thus,

$$\overline{R}_T = \mathbb{E}\!\left[\sum_{t=1}^{T} \ell_{I_t,t}\right] \;-\; \min_{i=1,\ldots N} \mathbb{E}\!\left[\sum_{t=1}^{T} \ell_{it}\right] \;\le\; \frac{\ln N}{\eta_T} \;+\; \frac{M^2 N}{2} \sum_{t=1}^{T} \eta_t$$

the only adaptation to be made is wrt $T$ (as $M$ is assumed to be known)

The optimal constant $\eta$ would be

s.t. $\dfrac{\ln N}{\eta} = \dfrac{M^2 N}{2} T \eta$, that is,

$\eta$ proportional to $\dfrac{1}{M} \sqrt{\dfrac{\ln N}{N T}}$

$\hookrightarrow$ try $\eta_t = \dfrac{\gamma}{M} \sqrt{\dfrac{\ln N}{N t}}$   where $\gamma$ is to be optimized.

The final bound is

$$M\sqrt{N\ln N}\left(\frac{\sqrt{T}}{\gamma} + \frac{\gamma}{2}\sum_{t=1}^{T}\frac{1}{\sqrt{t}}\right)$$

$$\underbrace{\qquad}_{\leq \int_{0}^{T}\frac{1}{\sqrt{t}}\,dt \;\leq\; 2\sqrt{T}}$$

$$\leq \sqrt{T}\left(\frac{1}{\gamma} + \gamma\right) = 2\sqrt{T} \quad \text{for } \gamma = 1$$

## Remarks / insights

* We heavily used above that the range $[0, M]$ is known

   • 0 to apply safely $e^{-x} \leq 1 - x + \frac{x^2}{2}$

   • $M$ to bound $\hat{\ell}_{j,t}^2 \leq M^2$

   • 0 and $M$ known because the $\eta_t$ are set based on the bounds obtained based on the $\eta_t$ previous two inequalities

* When the range $[m, M]$ is known, with $m \in \mathbb{R}$

   $\hookrightarrow$ Translate all losses by $-m$, eg, consider

   $$\hat{\ell}_{jt} = \frac{\ell_{I_t t} - m}{p_{jt}}\,\mathbb{1}_{\{I_t = j\}} \quad \text{and} \quad \eta_t = \frac{1}{M-m}\sqrt{\frac{\ln N}{Nt}}$$

   to get $\quad R_T \leq 2(M-m)\sqrt{TN\ln N}$

* What about an unknown range $[m, M]$? $\quad\rightsquigarrow$ See the homework #2!

## Additional elements $\quad\rightsquigarrow\quad$ See the extra lecture notes posted on the website

* With exponential weights, one can get high-probability bounds on the true regret, of the same order of magnitude:

   wp $1-\delta$, $\quad R_T \leq \square\,(M-m)\sqrt{TN\ln(N/\delta)}$

* These exists an algorithm (called INF) s.t. $R_T$ is controlled (in $\mathbb{E}$ or wp $1-\delta$) by $O(\sqrt{TN})$, ie, no $\sqrt{\ln N}$ term, against oblivious individual seq. $\ell_{jt}$