

Le premier exercice de l'examen vous fera prouver le résultat suivant :

Distribution-free (ie uniform) lower bounds.

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

$$\inf_{\text{strategies } \Psi} \sup_{\substack{x_1, \dots, x_K \\ \text{in } \mathcal{P}(\mathcal{E}_T)}} R_T \geq \frac{1}{20} \sqrt{T K}.$$

↑
and even:
 $\sup_{\substack{\text{over } x_1, \dots, x_K \\ \text{being Bernoulli distributions}}} R_T$

The MOSS strategy

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

Index policy relying on

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

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 $A_t \in \arg \max_{a=1 \dots K} U_a(t-1)$

Difference to UCB: we replace the exploration bonus

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

↳ no exploration after a was pulled sufficiently often (γ_K times)

We prove a distribution-free bound:

Theorem: MOSS is such that $\sup_{\substack{1 \leq t \leq T \\ \text{distributions} \\ \text{over } [0^T]}} R_T \leq K-1 + 45 \sqrt{KT}$

over $[0^T]$

(the constant 45 can be improved)

but indeed minimax optimal as its name indicates!

Proof.First step:

$$U_{\alpha^*}(t-i) \leq U_{A_t^*}(t-i) \quad \text{by definition of } A_t^* \text{ as an argmax}$$

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

$$\leq K-1 + \sum_{t=K+1}^T E[\mu^* - U_{\alpha^*}(t-i)] + \sum_{t=K+1}^T E[U_{A_t^*}(t-i) - \mu_{A_t^*}]$$

↑
if we played each arm once in the first K steps, and at most $K-1$ were suboptimal

$$\leq \sqrt{KT} + \sum_{t=K+1}^T E[(U_{\alpha^*}(t-i) - \mu_{A_t^*})^+]$$

Second step:Control of each $E[\mu^* - U_{\alpha^*}(t)]$ term by $20\sqrt{\frac{K}{t}}$ We write $E[\mu^* - U_{\alpha^*}(t)]$
for $t > K$

$$\leq E[(\mu^* - U_{\alpha^*}(t))^+]$$

$$\leq \sum_{l=0}^{+\infty} E[(\mu^* - U_{\alpha^*}(t))^+ \mathbb{1}_{\{N_{\alpha^*}(t) \in [x_{l+1}, x_l]\}}] \quad \text{where } x_l = \beta^{-l} t/K \text{ for some fixed } \beta > 1 \text{ and } l = 0, 1, 2, \dots$$

$$+ E[(\mu^* - U_{\alpha^*}(t))^+ \mathbb{1}_{\{N_{\alpha^*}(t) > t/K\}}]$$

$$\text{Now, } U_{\alpha^*}(t) = \hat{\mu}_{\alpha^*}(t) + \begin{cases} 0 & \text{if } N_{\alpha^*}(t) > t/K \\ \sqrt{\frac{1}{2N_{\alpha^*}(t)} \ln \left(\frac{t}{K\hat{\mu}_{\alpha^*}(t)} \right)} & \text{if } N_{\alpha^*}(t) \in [x_{l+1}, x_l] \end{cases}$$

denoted
by ε_l

Therefore,

$$(*) \quad E[\mu^* - U_{\alpha^*}(t)] \leq E[(\mu^* - \hat{\mu}_{\alpha^*}(t))^+ \mathbb{1}_{\{N_{\alpha^*}(t) > t/K\}}] + \sum_{l=0}^{+\infty} E[(\mu^* - \hat{\mu}_{\alpha^*}(t) - \varepsilon_l)^+ \mathbb{1}_{\{N_{\alpha^*}(t) \in [x_{l+1}, x_l]\}}]$$

$$\text{Lemma: } E[(\mu^* - \hat{\mu}_{\alpha^*}(t) - \varepsilon)^+ \mathbb{1}_{\{N_{\alpha^*}(t) \geq n_0\}}] \leq \frac{1}{\sqrt{n_0}} e^{-2n_0 \varepsilon^2}$$

Proof of the lemma:

$$E[(\mu^* - \hat{\mu}_{\alpha^*}(t) - \varepsilon)^+ \mathbb{1}_{\{N_{\alpha^*}(t) \geq n_0\}}] = \int_0^{+\infty} P\{(\mu^* - \hat{\mu}_{\alpha^*}(t) - \varepsilon) \geq u \text{ and } N_{\alpha^*}(t) \geq n_0\} du$$

$$= \int_0^{+\infty} \mathbb{P}\{ Z_{t+u}^* \geq (\varepsilon + u) N_{\alpha^*}(t) \text{ and } N_{\alpha^*}(t) \geq n_0 \} du$$

where $Z_t^* = N_{\alpha^*}(t) (\mu^* - \hat{\mu}_{\alpha^*}(t)) = \sum_{s=1}^t (\mu^* - \bar{y}_s) \mathbb{1}_{\{A_s = a^*\}}$
 is a wantingale,

and for all $x \in \mathbb{R}$, $S_{x,t} = e^{xZ_t^* - x^2/8 N_{\alpha^*}(t)}$

is a superwantingale.

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

$$= \int_0^{+\infty} \mathbb{P}\{ e^{xZ_t^* - x^2/8 N_{\alpha^*}(t)} \geq \exp(N_{\alpha^*}(t)(x(\varepsilon+u) - \frac{x^2}{8})) \text{ and } N_{\alpha^*}(t) \geq n_0 \} du$$

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

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

is independent of l , which
 will be useful in other
 proofs!

$$\begin{aligned} &\leq e^{-2n_0\varepsilon^2} \int_0^{+\infty} e^{-2n_0u^2} \mathbb{E}[S_{4(\varepsilon+u), t} \mathbb{1}_{\{N_{\alpha^*}(t) \geq n_0\}}] du \\ &\quad \text{where } \mathbb{E}[S_{4(\varepsilon+u), t}] \leq 1 \end{aligned}$$

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

$$\leq e^{-2n_0\varepsilon^2} \int_0^{+\infty} e^{-2n_0u^2} du = e^{-2n_0\varepsilon^2} \sqrt{\frac{\pi}{8n_0}}$$

integral of
 a Gaussian density,
 up to the normalization factor

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

Substituting the lemma in (*):

$$\mathbb{E}[\mu^* - U_{\alpha}(t)] \leq \sqrt{\frac{K}{t}} + \sum_{l=0}^{+\infty} \underbrace{\frac{1}{\sqrt{x_{\alpha l}}} e^{-2x_{\alpha l}} \frac{\epsilon_l^2}{e}}_{\underbrace{\frac{1}{\sqrt{x_{\alpha l}}} \exp(-2x_{\alpha l}) \frac{1}{2x_{\alpha l}} \ln \left(\frac{t}{Kx_{\alpha l}}\right)}_{\underbrace{\frac{1}{\beta}}_{\ln \beta = l \ln \beta}}}$$

where

$$\sum_{l \geq 0} \beta^{l(\frac{1}{2} - \frac{1}{\beta})} \text{ is } < \infty$$

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

$$= \sqrt{\frac{K}{t}} \beta^{(l+1)/2} \exp(-\frac{l}{\beta} \ln \beta)$$

$$= \sqrt{\frac{K}{t}} \beta^{1/2 + l(\frac{1}{2} - \frac{1}{\beta})}$$

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

$$\sum_{l=0}^{+\infty} \left(\frac{3}{2}\right)^{1/2 + l(\frac{1}{2} - \frac{1}{3})} = \sqrt{\frac{3}{2}} \sum_{l=0}^{+\infty} \alpha^{l+1} = \frac{1}{1-\alpha} \sqrt{\frac{3}{2}} \leq 9$$

where $\alpha = \left(\frac{3}{2}\right)^{\frac{1}{2} - \frac{1}{3}} \in (0, 1)$

All in all: we obtain a $\sqrt{\frac{K}{t}} + 9\sqrt{\frac{K}{t}} = 20\sqrt{\frac{K}{t}}$ bound, as claimed.

Third step:

$$\sum_{t=K+1}^T \mathbb{E}[(U_{A_{t+1}}(t-1) - \mu_{A_{t+1}} - \sqrt{\frac{K}{T}})^+] \leq 4\sqrt{KT}$$

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

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

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

We now use $(U_a(t) - \mu_a - \sqrt{\frac{K}{T}})^+ \leq (\hat{\mu}_a(t) - \mu_a - \sqrt{\frac{K}{T}})^+ + \int_0^\infty \mathbb{1}_{\{N_a(t) \geq t\}} \frac{1}{\sqrt{2\pi N_a(t)}} \ln \left(\frac{t}{K N_a(t)} \right)$

and get therefore the upper bound

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

$$+ \sum_{a=1}^K \left(\sum_{l=1}^{T-1} \sqrt{\frac{1}{2l}} \ln \left(\frac{T}{Kl} \right) \right)$$

$$\sqrt{\frac{1}{2N_a(t)} \ln \left(\frac{T}{KN_a(t)} \right)}$$

if $N_a(t) \leq \frac{T}{K}$
also smaller than

$$\mathbb{E}\left[\sum_{t=k}^{T-1} \mathbb{1}_{\{N_a(t)=l\}} \mathbb{1}_{\{A_{t+1}=a\}}\right]$$

We will repeatedly use that

$$\sum_{t=k}^{T-1} \mathbb{1}_{\{A_{t+1}=a_j\}} \mathbb{1}_{\{N_t(t)=l\}} \leq 1 \quad (\text{ie, disjoint union})$$

If $N_t(t)$ increases by 1 whenever a is played

$$\begin{aligned} \text{Also, } \sum_{l=1}^{\lceil T/k \rceil} \sqrt{\frac{1}{2k} \ln \left(\frac{T}{Kx} \right)} &\leq \int_0^{\lceil T/k \rceil} \sqrt{\frac{1}{2x} \ln \left(\frac{T}{Kx} \right)} dx \\ &\leq \sqrt{\frac{T}{2K}} \int_1^{+\infty} u^{-3/2} \sqrt{\ln u} du \\ &= \sqrt{\frac{T}{2K}} \int_0^1 u^{1/2} e^{-u/2} du \quad \text{by the change of variable } u = e^{v^2} \\ &= \sqrt{\pi} \sqrt{\frac{T}{K}} \end{aligned}$$

Summarizing what we proved so far:

$$\sum_{t=k}^{T-1} \mathbb{E}[(\hat{\mu}_a(t) - \mu_a - \sqrt{k_x})^+] \leq \sqrt{\pi} \sqrt{KT} + \underbrace{\sum_{a=1}^K \sum_{l=1}^{\lceil T/k \rceil} \sum_{t=k}^{T-1} \mathbb{E}[(\hat{\mu}_a(t) - \mu_a - \sqrt{k_x})^+]}_{\leq \sqrt{\pi/2} \sqrt{T_k} \text{ for each } a} \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_t(t)=l\}}$$

We resort again to $Z_{ait} = N_t(t)(\hat{\mu}_a(t) - \mu_a)$

and

$$S_{x,t}^{(a)} = e^{xz_{ait}} - \frac{x^2}{8} N_t(t) \quad \text{super martingale}$$

where $x = 4(\sqrt{k_x} + u)$

For each a_j ,

$$\begin{aligned} \sum_{l=1}^{\lceil T/k \rceil} \sum_{t=k}^{T-1} \mathbb{E}[(\hat{\mu}_a(t) - \mu_a - \sqrt{k_x})^+ \mathbb{1}_{\{A_{t+1}=a_j\}} \mathbb{1}_{\{N_t(t)=l\}}] \\ = \sum_{l=1}^{\lceil T/k \rceil} \sum_{t=k}^{T-1} \int_0^{+\infty} \Pr[x Z_{ait} \geq N_t(t) \left(x(u + \sqrt{k_x}) + A_{t+1} = a_j \wedge N_t(t) = l \right)] du \\ \stackrel{Markov-Chebyshev}{\leq} \sum_{l=1}^{\lceil T/k \rceil} \sum_{t=k}^{T-1} \int_0^{+\infty} e^{-du} \underbrace{e^{-\frac{d}{2} l(l(u + \sqrt{k_x}))^2}}_{\leq e^{-2lu^2} e^{-2lK_x}} \underbrace{\mathbb{E}[S_{x,t}^{(a_j)} \mathbb{1}_{\{A_{t+1}=a_j\}} \mathbb{1}_{\{N_t(t)=l\}}]}_{\text{Be sum over } t \text{ of these will be } \leq 1} du \\ \text{issue: this depends on } t \dots \\ \text{but can be replaced in some sense by } S_{x,0}^{(a_j)} = 1 \end{aligned}$$

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

$$\mathbb{P}\left\{\sup_{t \geq 0} S_{\leq t}^{(a)} \geq c\right\} \leq \frac{E[S_{\leq 0}^{(a)}]}{c} = \frac{1}{c}$$

see also an alternative treatment on the next page

Then,

$$\sum_{l=1}^T \sum_{t=k}^{T-1} \mathbb{E}[(\hat{\mu}_k(t) - \mu_k - \sqrt{k_T})^+ \mathbf{1}_{A_{t,k}=a} \mathbf{1}_{N_k(t)=l}]$$

$$\text{as before!} \quad = \sum_{l=1}^{+\infty} \sum_{t=k}^{T-1} \int_0^{+\infty} \mathbb{P}\left\{S_{\leq t}^{(a)} \geq e^{2l(u + \sqrt{k_T})^2} \text{ and } A_{t,k}=a \text{ and } N_k(t)=l\right\} du$$

$$\leq \sum_{l=1}^T \int_0^{+\infty} \sum_{t=k}^{T-1} \mathbb{P}\left(\sup_{t \geq 0} S_{\leq t}^{(a)} \geq e^{2l(u + \sqrt{k_T})^2} \text{ and } A_{t,k}=a \text{ and } N_k(t)=l\right) du$$

$$\leq \sum_{l=1}^T \int_0^{+\infty} \mathbb{P}\left(\sup_{t \geq 0} S_{\leq t}^{(a)} \geq e^{2l(u + \sqrt{k_T})^2}\right) du$$

f. disjoint which!

$$\text{Doob's maximal inequality} \leq \sum_{l=1}^T \int_0^{+\infty} e^{-2l(u + \sqrt{k_T})^2} du \leq \sum_{l=1}^{+\infty} \frac{1}{\sqrt{e}} e^{-2lK/T}$$

$\leq e^{-2lu^2} \times e^{-2lK/T}$
and same treatment as
in the lemma of the first part
of the proof

This step is concluded by calculations :

$$\begin{aligned} \sum_{l=1}^T \frac{1}{\sqrt{e}} e^{-2lK/T} &\leq \int_0^T \frac{1}{\sqrt{e}} e^{-2xK/T} dx \\ &= \sqrt{\frac{1}{2K}} \int_0^T \frac{e^{-u}}{\sqrt{u}} du = \sqrt{\frac{T}{2K}} \int_0^{+\infty} e^{-v^2} dv = \sqrt{\frac{\pi}{2}} \sqrt{\frac{T}{K}} \end{aligned}$$

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

General conclusion : Final bound given by

$$\begin{aligned} K-1 + \left(\sum_{t=k+1}^T 2\sqrt{\frac{K}{t}} \right) + \sqrt{KT} + 4\sqrt{KT} &\leq K-1 + 5\sqrt{KT} + 2\sqrt{\int_0^T \sqrt{\frac{K}{t}} dt} \\ &= K-1 + 4\sqrt{KT} \end{aligned}$$

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

We were stuck at

$$\sum_{l=1}^{\infty} \sum_{t=k}^{T-1} \int_0^{+\infty} e^{-2lu^2} e^{-2lKt} E \left[S_{x_t}^{(a)} \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=l\}} \right] du \\ = \sum_{l=1}^{\infty} \int_0^{+\infty} e^{-2lu^2} e^{-2lKt} E \left[\sum_{t=k}^{T-1} S_{x_t}^{(a)} \mathbb{1}_{\{A_{t+1}=a\}} \mathbb{1}_{\{N_a(t)=l\}} \right] du \\ = S_{x_T T_l}^{(a)}$$

where T_l is given by:

$$T_l = \inf \{ t \in \mathbb{Q} | ... T \} : A_{t+1} = a \text{ and } N_a(t) = l \}$$

We should get $E[S_{x_T T_l}^{(a)}] \leq E[S_{x_0}^{(a)}] = 1$

from the optional stopping theorem (\Leftarrow théorème d'arrêt de Doob \Rightarrow) provided some verifications.

(T_l should be a bounded stopping time.)

Adversarial bandits.

(Rather stated in terms of losses than rewards!)

Setting:

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

1. The opponent and the decision-maker simultaneously choose $\ell_t = (\ell_{jt})_{j=1..N}$ and $I_t \sim p_t$, where $p_t \in \mathcal{P}(\{1, \dots, N\})$

2. The opponent gets to see p_t and I_t ; the decision-maker only observes $\ell_{I_t, t}$ (her own loss).

Regret:

$$R_T = \sum_{t=1}^T \ell_{I_t, t} - \min_{j=1..N} \sum_{t=1}^T \ell_{jt}$$

vs. Pseudo-regret: $\bar{R}_T = \mathbb{E}\left[\sum_{t=1}^T \ell_{I_t, t}\right] - \min_{j=1..N} \mathbb{E}\left[\sum_{t=1}^T \ell_{jt}\right]$

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

↑
Why \mathbb{E} ?
(p. ℓ_{jt} are
random variables, as
they depend on the first,
and in particular on I_t ,
 \dots, I_{t-1})

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

We actually rather shoot for high-probability bounds on R_T , but studying \bar{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 \bar{R}_T via exponential weights.

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

$$\hat{l}_{jt} = \frac{l_{I_t t}}{p_{jt}} \mathbb{1}_{\{I_t = j\}}$$

if $p_{jt} > 0$
(which we will assume)

auxiliary
randomizations
of opponent +
decision-maker

They are (conditionally) unbiased: denoting by $\mathcal{F}_{t-1} = \sigma(l_1, \dots, l_{t-1}, U_{t-1}, U_t, p_1, \dots, p_{t-1})$

The total information available at the beginning of round t

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

we have:

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

so that

$$\mathbb{E}[\hat{l}_{jt} | \mathcal{F}_{t-1}] = \frac{\hat{l}_{jt}}{p_{jt}} \mathbb{E}\left[\underbrace{\mathbb{1}_{\{I_t = j\}}}_{= p_{jt}} | \mathcal{F}_{t-1}\right] = \frac{\hat{l}_{jt}}{p_{jt}} p_{jt} = \hat{l}_{jt}$$

Since we assumed $p_{jt} > 0$...

Algorithm:

$p_1 = (\frac{1}{N}, \dots, \frac{1}{N})$ and for $t \geq 2$, $p_t = (p_{jt})_{j=1,\dots,N}$ is

defined as

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

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

↳ ensures indeed that $p_{jt} > 0$.

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

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

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

$$\bar{R}_T = \mathbb{E}\left[\sum_{t=1}^T l_{I_t t}\right] - \min_{i=1,\dots,N} \mathbb{E}\left[\sum_{t=1}^T l_{it}\right] \leq 2M \sqrt{T \ln N}$$

The proof is based on the following lemma.

Lemma: The exponentially weighted average strategy on losses

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

$$\text{with } \eta_t \downarrow, \quad \tilde{p}_{jt} = \exp(-\eta_t \sum_{s=1}^{t-1} \tilde{l}_{js}) / \sum_{k=1}^N \exp(-\eta_t \sum_{s=1}^{t-1} \tilde{l}_{ks}),$$

is such that

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

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

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

$$\text{where } \tilde{s}_t = \sum_{j=1}^N \tilde{p}_{jt} \tilde{l}_{jt} + \frac{1}{\eta_t} \ln \sum_{j=1}^N \tilde{p}_{jt} e^{-\eta_t \tilde{l}_{jt}}$$

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

$$\text{so that } \ln \sum_j \tilde{p}_{jt} e^{-\eta_t \tilde{l}_{jt}} \leq \ln (1 - \eta_t \sum_j \tilde{p}_{jt} \tilde{l}_{jt} + \frac{\eta_t^2}{2} \sum_j \tilde{p}_{jt} \tilde{l}_{jt}^2)$$

$$\leq -\eta_t \sum_j \tilde{p}_{jt} \tilde{l}_{jt} + \frac{\eta_t^2}{2} \sum_j \tilde{p}_{jt} \tilde{l}_{jt}^2$$

$\ln(1+x) \leq x$ for $x > -1$

hence the stated bound.

Proof (of the theorem): We have no control on how large the \tilde{l}_{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 $\tilde{l}_{jt} \in [0, M_T]$... 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 η_t can be picked as $\ln N / \sum_{s=1}^{t-1} \delta_s$ or in terms of the upper bounds on the δ_s (we choose the latter version for the sake of conciseness).
 \hookrightarrow see below!

The lemma yields for the \hat{l}_{it} :

$$\begin{aligned} \sum_{t,j} p_{jt} \hat{l}_{jt} - \min_{i=1, \dots, N} \sum_{t=1}^T \hat{l}_{it} &\leq \frac{\ln N}{\eta_T} + \sum_{t=1}^T \frac{\eta_t}{2} \sum_j p_{jt} \hat{l}_{jt}^2 \\ &\quad \downarrow \text{by definition of the } \hat{l}_{it} \\ &= \sum_j l_{I_t t} \frac{p_{jt}}{p_{I_t t}} \mathbb{1}_{\{I_t = j\}} = l_{I_t t} \\ &\quad \downarrow \text{similar treatment:} \\ &= \sum_j l_{I_t t}^2 \frac{1}{p_{I_t t}} \mathbb{1}_{\{I_t = j\}} \\ &\leq M^2 \sum_j \frac{\mathbb{1}_{\{I_t = j\}}}{p_{I_t t}} \end{aligned}$$

To simplify even further the choice of the η_t , we first take E of both sides:

$$\begin{aligned} E\left[\sum_t l_{I_t t}\right] - E\left[\min_i \sum_t \hat{l}_{it}\right] &\leq \frac{\ln N}{\eta_T} + \frac{M^2}{2} \sum_{t=1}^T \eta_t \sum_{j=1}^N E\left[\frac{\mathbb{1}_{\{I_t = j\}}}{p_{I_t t}}\right] \\ &\leq \min_i \sum_{t=1}^T E[\hat{l}_{it}] \\ &= E[\hat{l}_t] \quad \text{by the tower rule and the fact that } \hat{l}_t \text{ is conditionally unbiased.} \end{aligned}$$

Thus,

$$\bar{R}_T = E\left[\sum_{t=1}^T l_{I_t t}\right] - \min_{i=1, \dots, N} E\left[\sum_{t=1}^T \hat{l}_{it}\right] \leq \frac{\ln N}{\eta_T} + \frac{M^2 N}{2} \sum_{t=1}^T \eta_t$$

$\underbrace{\quad}_{\text{the only adaptation to be made is w.r.t } T} \quad \underbrace{\quad}_{\text{(as } M \text{ is assumed to be known)}}$

The optimal constant η would be

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

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

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

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

$$\leq \int_0^T \frac{1}{\sqrt{t}} dt \leq 2\sqrt{T}$$

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

Remarks / insights:

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

M to bound $\ell_{I_t, t}^2 \leq M^2$
and to pick η_t

→ The range $[0, M]$ needs to be known; it could be a general range $[m, M]$, in which case, we would translate all losses by $-m$,
e.g. consider $\hat{\ell}_{I_t, t} = \frac{\ell_{I_t, t} - m}{p_{I_t}} \mathbf{1}_{I_t = j}$

and $\eta_t = \frac{1}{M-m} \sqrt{\frac{\ln N}{Nt}}$

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

* Last-minute edit:
I think I have a solution
for R_T based on parts of
the additional document

↳ Adaptation to the range in adversarial bandits is
actually a difficult issue (I would need to
think more about it...)*

↳ extra points for whoever finds it!
→ In the next page you'll get a brief view of how to get high-probability bounds on
the true regret, of the same order of magnitude:

w.p. 1-δ, $R_T \leq \square (M-m) \sqrt{T N \ln(N/\delta)}$

→ There exists an algorithm (INF) s.t. R_T is controlled (in E or with
high prob.) by $O(\sqrt{N})$ against individual sequences ℓ_{I_t} [ie: sequences
fixed in advance, ↳ not chosen by an opponent who reacts] ↳
ie: no $\sqrt{N \ln N}$ needed ↳ « oblivious » opponent ↳ The ℓ_{I_t} are not
random variables in this case