

Pour aller plus loin sur le thème des "adversarial bandits" (à défaut d'avoir pu mieux voir ce sujet en cours; ce qui suit ne sera pas testé à l'examen)



- Distribution-free regret bounds for stochastic bandits
- Adversarial bandits

To that end we restrict our attention to the model $\mathcal{D} = \mathcal{P}(\mathcal{Q}_1)$, the set of all probability distributions over \mathcal{Q}_1 .

Stochastic bandits

With each arm a is associated $\bar{g}_a \in \mathcal{P}(\mathcal{Q}_1)$

For $t=1, 2, \dots$

- The decision maker picks $I_t \in \{1, \dots, K\}$
- Her reward \bar{y}_t , which is such that $\bar{y}_t | I_t \sim \bar{g}_{I_t}$, is her only piece of information

Aim: control the regret

$$\bar{R}_T = \bar{T} \max_{a=1..K} E(\bar{g}_a) - E\left[\sum_{t=1}^T \bar{y}_t\right]$$

Adversarial bandits

An opponent selects the payoffs g_{jt}

For $t=1, 2, \dots$

- The opponent picks $(g_{jt} - g_{kt}) \in \mathcal{Q}_1^K$ while, simultaneously,
- The decision-maker picks $I_t \in \{1, \dots, K\}$
- Her payoff is $g_{I_t t}$ and this is the only piece of information she gets

Aim: control the regret

$$R_T = \max_{k=1..K} \sum_{t=1}^T g_{kt} - \sum_{t=1}^T g_{I_t t}$$

Typical adversarial results

(Auer, Cesa-Bianchi, Freund, Schapire, 2002, later improved by Audibert and Bubeck, 2009):

Strategies

such that for all opponents picking gains in \mathcal{Q}_1 ,

{ with probability at least $1-\delta$,

$$E[R_T] \leq C \sqrt{T \ln K}$$

where

the probability and E are w.r.t. decision-maker's internal randomization

for all $T \geq 1$,

$$R_T \leq C \sqrt{T K \ln(K/\delta)}$$

for some numerical constant C

For "oblivious" opponents (i.e., when the g_{jt} do not "react" to the decision-maker's actions): the $\sqrt{\ln K}$ can be dropped.

It is in particular

the case when $g_{jt} \sim \bar{g}_k$ $\forall t$ in an independent way

In this stochastic model:

$$\begin{aligned} E[R_T] &= E\left[\max_{k=1..K} \sum_{t=1}^T g_{kt}\right] - E\left[\sum_{t=1}^T g_{It,t}\right] \quad \text{→ } g_{It,t} \text{ is } y_t \\ &> T \max_{k=1..K} E[g_{k,1}] - E\left[\sum_{t=1}^T y_t\right] \\ &= T \max_{k=1..K} E(\bar{g}_k) - E\left[\sum_{t=1}^T y_t\right] = \bar{R}_T \end{aligned}$$

The adversarial results entail in particular that there exists a strategy of the decision-maker such that

$$\sup_{\vec{y}_1, \dots, \vec{y}_K \in \mathcal{P}([0,1])} \bar{R}_T \leq \sup_{\substack{\text{opponents} \\ \text{picking } g_{It} \in [0,1]}} E[R_T] \leq C \sqrt{T K} \quad \text{for some numerical constant } C$$

while the ranking of lower bounds is:

For all (randomized) strategies of the decision-maker,
for all $K \geq 2$ and $T \geq K/5$,

$$\sup_{\substack{\text{opponents} \\ \text{picking } g_{It} \in [0,1]}} E[R_T] \geq \sup_{\substack{\text{individual} \\ \text{sequences } g_{It} \in [0,1]}} E[R_T] \geq \sup_{\vec{y}_1, \dots, \vec{y}_K \in \mathcal{P}([0,1])} \bar{R}_T \geq \frac{1}{20} \sqrt{T K}$$

↑
and even:
 $\sup_{\substack{\text{over } \vec{y}_1, \dots, \vec{y}_K \\ \text{being Bernoulli} \\ \text{distributions}}} E[R_T]$
as proved
last week

Conclusion:

The minimax rates of the regret for stochastic bandits on $[0,1]$ or against oblivious opponents picking rewards in $[0,1]$ are \sqrt{TK} .

But: What is the minimax rate against general recursive opponents?

- Should/Can the $\sqrt{TK/K}$ upper bound be improved?
- Should/Can the $\sup_{\text{opponents}} E[R_T]$ lower bound be improved?

↓
look for sequences of payoffs g_{kt} with real and strong correlations/dependencies in the past

Adversarial bandits:high-probability regret bounds-(A brief and
suboptimal
view...)Setting (reminder): at each round $t=1, 2, \dots$

1. the opponent and the decision-maker simultaneously choose $l_t = (l_{jt})_{j=1..N}$ with $l_{jt} \in [0, M]$ and $I_t \sim p_t$, where $p_t \in \mathcal{P}_{\{1..N\}^p}$;
2. the opponent gets to see p_t and I_t ; the decision-maker only observes $l_{I_t, t}$ (her own loss).

↳ She wants to control her regret

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

She resorts to the (conditionally) unbiased estimators

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

Denoting by $\mathcal{F}_{t-1} = \sigma(p_1, I_1, l_1, \dots, t-1)$
 we showed: and $\mathcal{F}_0 = \{\emptyset, \Omega\}$

$$\forall t \geq 1, \quad \forall j \in \{1, \dots, N\}, \quad E[\hat{l}_{jt} | \mathcal{F}_{t-1}] = l_{jt}$$

Main technical ingredient needed in the proof:Be able to relate $\sum_{t=1}^T \hat{l}_{jt}$ to $\sum_{t=1}^T l_{jt}$ with high probability.

We will use martingale inequalities but will need a control from

below on the p_t . Hence the Exp3 algorithm:for $t \geq 2$, $\forall j$,

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

↑
exploitation term ↑
exploration term

Hence the name Exp₃ for: exponential weights for exploration and exploitation

Then: $\hat{l}_{jt} \in [0, \frac{MN}{\gamma_t}]$ as $p_{jt} \geq \gamma_t/N$ and $l_{jt} \in [0, M]$

The Hoeffding-Azuma inequality ensures that with probability at least $1-\delta$,

$$(*) \quad \sum_{t=1}^T \hat{l}_{jt} \leq \sum_{t=1}^T \underbrace{\mathbb{E}[\hat{l}_{jt} | \mathcal{F}_{t-1}]}_{= \bar{l}_{jt}} + MN \sqrt{\sum_{t=1}^T \frac{1}{2\gamma_t^2} \ln \frac{1}{\delta}}$$

depends on
the range

A sharper bound is in terms of the conditional variances:

$$\text{Var}_{\mathcal{F}_{t-1}}(\hat{l}_{jt}) \leq \mathbb{E}[\hat{l}_{jt}^2 | \mathcal{F}_{t-1}] \leq M^2 \mathbb{E}\left[\frac{1_{\{l_{jt} > 0\}}}{p_{jt}^2} | \mathcal{F}_{t-1}\right]$$

↑
as $l_{jt}^2 \leq M^2$ = $M^2 p_{jt}$ $\leq \frac{M^2 N}{\gamma_t}$

↳ Bernstein's inequality for martingales ensures that

with probability at least $1-\delta$,

$$(**) \quad \sum_{t=1}^T \hat{l}_{jt} \leq \sum_{t=1}^T \bar{l}_{jt} + O\left(M\sqrt{N} \sqrt{\sum_{t=1}^T \frac{1}{\gamma_t} \ln \frac{1}{\delta}}\right)$$

↑

→ With (*) we would get a regret bound of order $T^{3/4}$

We will gain orders of magnitude as γ_t will be of the form $t^{-\alpha}$

while with (**) we can get the desired \sqrt{T} , after an additional twist (we will first exhibit a $T^{2/3}$ bound and only hint at the \sqrt{T} rate as an exercise).

Bernstein's inequality for martingale.

We proved Bernstein's lemma with true expectations but given its proof (where we used only the monotonicity of E), it appears that we can replace all E by $E[\cdot | \mathcal{G}_j]$.

But I realized meanwhile that my proof was suboptimal as it used a lower bound m on the random variable X at hand, which is an inconvenient assumption. Let's re-do it.

Lemma: Let X be a random variable with $X - E[X|G_j] \leq M$ and G_j a σ -algebra; then :

$$\forall \eta \geq 0, \quad \ln E[e^{\eta X} | G_j] \leq \eta E[X|G_j] + \frac{e^{\eta M} - \eta M - 1}{M^2} \text{Var}_{G_j}(X)$$

or put differently,

$$E[e^{\eta(X - E[X|G_j])} | G_j] \leq \exp\left(\frac{e^{\eta M} - \eta M - 1}{M^2} \text{Var}_{G_j}(X)\right) \quad \uparrow \text{conditional variance of } X$$

Proof: $\eta(X - E[X|G_j]) \leq \eta M$ and $x \in \mathbb{R} \mapsto \frac{e^x - x - 1}{x^2}$ is increasing,

$$\text{so that } e^{\eta(X - E[X|G_j])} - \eta(X - E[X|G_j]) - 1 \leq \frac{e^{\eta M} - \eta M - 1}{M^2} \frac{e^{\eta M} - \eta M - 1}{M^2}$$

Taking $E[\cdot | G_j]$:

$$E[e^{\eta(X - E[X|G_j])} | G_j] - 1 \leq \underbrace{\text{Var}_{G_j}(X)}_{\ln u \leq u-1} \frac{e^{\eta M} - \eta M - 1}{M^2}$$

$$\hookrightarrow \ln E[e^{\eta(X - E[X|G_j])} | G_j] \geq \ln \frac{e^{\eta M} - \eta M - 1}{M^2}, \quad \text{which concludes the proof.}$$

Theorem: Let $(\mathcal{F}_t)_{t \geq 0}$ be a filtration and let $(X_t)_{t \geq 1}$ be a sequence of adapted random variables, with $X_t - E[X_t | \mathcal{F}_{t-1}] \leq M$ a.s., $\forall t$.

— probabilistic version: $\forall \varepsilon > 0, \quad \forall V > 0,$

$$P\left\{ \sum_{t=1}^T X_t - \sum_{t=1}^T E[X_t | \mathcal{F}_{t-1}] \geq \varepsilon \quad \text{and} \quad \sum_{t=1}^T \text{Var}_{\mathcal{F}_{t-1}}(X_t) \leq V \right\} \leq \exp\left(-\frac{\varepsilon^2}{2V + \frac{2}{3}M\varepsilon}\right)$$

- Statistical version :

$$\mathbb{P} \left\{ \sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}] \geq \sqrt{2V \ln \frac{1}{\delta}} + \frac{2}{3} M \ln \frac{1}{\delta} \right\} \leq \delta$$

where V is a numerical constant > 0

Proof: $\eta > 0$ $M_t = \exp \left(\eta \sum_{s=1}^t (X_s - \mathbb{E}[X_s | \mathcal{F}_{s-1}]) - (e^{\eta M} - e^{\eta M-1})/M^2 \sum_{s=1}^t \text{Var}_{\mathcal{F}_{s-1}}(X_s) \right)$

by Bernstein's lemma, $(M_t)_t$ is an $(\mathcal{F}_t)_t$ -adapted supermartingale.

Thus $\mathbb{P} \left\{ \underbrace{\sum_{t=1}^T (X_t - \mathbb{E}[X_t | \mathcal{F}_{t-1}])}_{\text{not. } S_T} \geq \varepsilon \text{ and } \underbrace{\sum_{t=1}^T \text{Var}_{\mathcal{F}_{t-1}}(X_t)}_{\text{not. } \sigma_T^2} \leq V \right\}$

$$\begin{aligned} &= \mathbb{P} \left\{ \underbrace{e^{\eta S_T} - (e^{\eta M} - e^{\eta M-1})/M^2 \sigma_T^2}_{= M_T} \geq e^{\eta \varepsilon - (e^{\eta M} - e^{\eta M-1})/M^2 \sigma_T^2} \text{ and } \sigma_T^2 \leq V \right\} \\ &\leq \mathbb{P} \left\{ M_T \geq \exp(\eta \varepsilon - (e^{\eta M} - e^{\eta M-1})/M^2 V) \right\} \end{aligned}$$

\leq Markov's Ineq. $\underbrace{\exp(-\eta \varepsilon + (e^{\eta M} - e^{\eta M-1})/M^2 V)}_{\text{to be optimized over } \eta > 0} \underbrace{\mathbb{E}[M_T]}_{\leq \mathbb{E} M_0 = 1 \text{ as } (\eta_t)_t \text{ is a supermartingale.}}$

$$f(x) = -x\varepsilon + \frac{V}{M} (e^x - x - 1) \quad \text{where } \eta M = x$$

$$f'(x) = -\varepsilon + \frac{V}{M} (e^x - 1)$$

$$f''(x) = \frac{Ve^x}{M} > 0 \quad \text{unique minimizer } x^* \text{ s.t.}$$

$$f'(x^*) = 0, \text{ i.e., } e^{x^*} = 1 + \frac{M\varepsilon}{V}$$

$$x^* = \ln \left(1 + \frac{M\varepsilon}{V} \right)$$

$$\hookrightarrow \eta^* = \frac{1}{M} \ln \left(1 + \frac{M\varepsilon}{V} \right)$$

The bound is

$$\begin{aligned} &\exp \left(-\frac{\varepsilon}{M} \ln \left(1 + \frac{M\varepsilon}{V} \right) + \frac{V}{M^2} \left(\left(1 + \frac{M\varepsilon}{V} \right) - \ln \left(1 + \frac{M\varepsilon}{V} \right) - 1 \right) \right) \\ &= \exp \left(\frac{\varepsilon}{M} - \frac{\varepsilon}{M} \ln \left(1 + \frac{M\varepsilon}{V} \right) - \frac{V}{M^2} \ln \left(1 + \frac{M\varepsilon}{V} \right) \right) \end{aligned}$$

$$= \exp\left(-\frac{V}{M^2} \left(-\frac{M\varepsilon}{V} + \left(\frac{M\varepsilon}{V} + 1\right) \ln\left(1 + \frac{M\varepsilon}{V}\right)\right)\right)$$

$$= \exp\left(-\frac{V}{M^2} h\left(\frac{M\varepsilon}{V}\right)\right)$$

where $h(u) = (1+u) \ln(1+u) - u$

We conclude by noting that $\forall u \in \mathbb{R}$, $h(u) \geq \frac{u^2}{2 + \frac{2}{3}u}$:

$$\frac{V}{M^2} h\left(\frac{M\varepsilon}{V}\right) \geq \frac{V}{M^2} \left(\frac{M\varepsilon}{V}\right)^2 / \left(2 + \frac{2}{3} \frac{M\varepsilon}{V}\right)$$

$$= \frac{\varepsilon^2}{V} / \left(2 + \frac{2}{3} \frac{M\varepsilon}{V}\right).$$

\uparrow
probabilistic version

Statistical version

\downarrow It suffices to show that for $\varepsilon = \sqrt{2V \ln \frac{1}{\delta}} + \frac{2}{3} M \ln \frac{1}{\delta}$,

we have $\exp(-\varepsilon^2 / (2V + \frac{2}{3} M\varepsilon)) \leq \delta$.

Indeed,

$$\begin{aligned} \varepsilon^2 &= \left(\sqrt{2V \ln \frac{1}{\delta}} + \frac{2}{3} M \ln \frac{1}{\delta}\right) \varepsilon \quad \text{where } \varepsilon \geq \sqrt{2V \ln \frac{1}{\delta}} \\ &\geq 2V \ln \frac{1}{\delta} + \frac{2}{3} M \varepsilon \ln \frac{1}{\delta} \\ &= \left(2V + \frac{2}{3} M \varepsilon\right) \ln \frac{1}{\delta}. \end{aligned}$$

Application: Performance bound for Exp3

Theorem. For well-chosen sequences of $\eta_t \downarrow$ and $\gamma_t \downarrow$, we have, with probability $\geq 1-\delta$,

$$R_T = \sum_{t=1}^T l_{I_t t} - \min_{i=1..N} \sum_{t=1}^T \hat{l}_{it} \leq O(T^{2/3} \ln \frac{N}{\delta})$$

Remarks: - Not quite the \sqrt{T} rate we wanted! ↳ see Exercise to see how to correct this.

- It would be even worse with Hoeffding-Azuma ($T^{3/4}$ rate).

- Main issue: the deviation term $\sum \hat{l}_{it} - \sum l_{it} \leq \dots \rightarrow$ to be improved.

Proof: Let $\tilde{p}_{jt} = \exp(-\eta_t \sum_{s=1}^{t-1} \hat{l}_{js}) / \sum_{k=1}^N \exp(-\eta_t \sum_{s=1}^{t-1} \hat{l}_{ks})$

(so that $p_{jt} = (1-\gamma_t) \tilde{p}_{jt} + \gamma_t / N$).

A lemma used already in the proof of the expected bound shows

that

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

thus $\sum_{t=1}^T \sum_{j=1}^N \left(\underbrace{(1-\gamma_t) \tilde{p}_{jt} + \frac{\gamma_t}{N}}_{\leq 1} \right) \hat{l}_{jt} - \min_{i=1..N} \sum_{t=1}^T \hat{l}_{it}$

$$\leq \sum_{t=1}^T \frac{\gamma_t}{N} \sum_j \hat{l}_{jt} + \left(\frac{\ln N}{\eta_T} + \sum_{t=1}^T \frac{\eta_t}{2} \sum_j \tilde{p}_{jt} \hat{l}_{jt}^2 \right)$$

↑ $(1-\gamma_t) \tilde{p}_{jt} \leq p_{jt}$

Therefore,

$$\begin{aligned} & \sum_{t=1}^T \sum_{j=1}^N p_{jt} \hat{l}_{jt} - \min_{i=1..N} \sum_{t=1}^T \hat{l}_{it} \\ & \leq \frac{\ln N}{\eta_T} + \frac{1}{2} \sum_{t=1}^T \frac{\eta_t}{1-\gamma_t} \sum_j p_{jt} \hat{l}_{jt}^2 + \sum_{t=1}^T \frac{\gamma_t}{N} \sum_j \hat{l}_{jt} \end{aligned}$$

We already saw that:

$$\sum_j p_{jt} \hat{l}_{jt} = l_{I_t t} \quad \text{and that} \quad \sum_j p_{jt} \hat{l}_{jt}^2 \leq M^2 \frac{\mathbb{1}_{I_t=j}}{p_{jt}}$$

Bernstein's inequality applied $3N$ times:

$$\boxed{\begin{aligned} \forall i, \sum_{t=1}^T \hat{l}_{it} &\leq \sum_{t=1}^T l_{it} + M\sqrt{N} \sqrt{\sum_{t=1}^T \frac{\eta_t}{\gamma_t} \ln \frac{3N}{\delta}} + \frac{2}{3} \frac{MN}{\gamma_T} \ln \frac{3N}{\delta} \\ \forall j, \sum_{t=1}^T \frac{\eta_t}{1-\gamma_t} - \frac{1}{Pj} &\leq \sum_{t=1}^T \frac{\eta_t}{1-\gamma_t} + \sqrt{2 \sum_{t=1}^T \frac{\eta_t^2}{(\gamma_t)^2} \frac{N}{\gamma_t} \ln \frac{3N}{\delta}} + \frac{2}{3} \max_{t \leq T} \frac{\eta_t}{1-\gamma_t} N \times \ln \frac{3N}{\delta} \\ &\text{conditional variance} \leq \frac{\eta_t}{Pj} \leq \frac{N}{\gamma_t} \\ \forall j, \sum_{t=1}^T \gamma_t \hat{l}_{jt} &\leq \sum_{t=1}^T \gamma_t l_{jt} + \sqrt{2 \sum_{t=1}^T \gamma_t M^2 N \ln \frac{3N}{\delta}} + \frac{2}{3} MN \ln \frac{3N}{\delta} \\ &\text{conditional variance} \leq \gamma_t^2 M^2 N / \gamma_t = M^2 N \gamma_t \end{aligned}}$$

↑ All inequalities holding at the same time with probability at least $1-\delta$
 (by the union bound).

The regret bound is of the form: with probability $\geq 1-\delta$,

$$\begin{aligned} R_T &= \sum_{t=1}^T l_{i^*_t, t} - \min_{i=1, \dots, N} \sum_{t=1}^T l_{it} \\ &\leq \frac{\ln N}{\gamma_T} + M\sqrt{N} \sqrt{\sum_{t=1}^T \frac{2}{\gamma_t} \ln \frac{3N}{\delta}} + M \sum_{t=1}^T \gamma_t + \frac{M^2}{2} \sum_{t=1}^T \frac{\eta_t}{1-\gamma_t} \\ &\quad [+ \text{MANY OTHER TERMS}] \end{aligned}$$

but looking only at the γ_t we see the issue:

$$\gamma_t \sim t^{-\alpha} \rightarrow (\sum \gamma_t)^{1/2} = O(t^{(-\alpha+1)/2})$$

$$\sum \gamma_t = O(t^{\alpha})$$

Choose α s.t. $-\alpha+1 = 2(\alpha+1)$ i.e., $\alpha = -1/3$

and get a $T^{2/3}$ rate...

↳ You can show that indeed, the $T^{2/3} \ln \frac{N}{\delta}$ rate is achievable here.

Exercise \sqrt{T} high-probability bound on the regret of Exp3.

Trick: bias the estimators!
(and translate)

$$\hat{l}_{jt} = M - \frac{M - l_{jt}}{p_{jt}} \mathbb{1}_{\{j \neq j^*\}} - \frac{\beta_t}{p_{jt}}$$

Algorithm:

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

Prove that for well-chosen η_t , γ_t and β_t :

$$R_T \leq O(\sqrt{T N \ln \frac{N}{\delta}}) \text{ with probability at least } 1 - \delta.$$

NOTE:

I'm sorry, I will not have time to write up corrections!