

The Hoeffding-Azuma inequality

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 (ie, $\forall t \geq 1$, X_t is \mathcal{F}_{t-1} -measurable), that are bounded: $\forall t, a_t \leq X_t \leq b_t$ a.s., where $a_t, b_t \in \mathbb{R}$.

Then (\approx probabilistic version)

$$\forall \varepsilon > 0, \quad \mathbb{P}\left\{\sum_{t=1}^T X_t - \sum_{t=1}^T \mathbb{E}[X_t | \mathcal{F}_{t-1}] \geq \varepsilon\right\} \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{t=1}^T (b_t - a_t)^2}\right)$$

or (\approx statistical version), totally equivalent

$\forall \delta \in (0, 1)$, with probability at least $1 - \delta$,

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

Note: Hoeffding's inequality is the special case when all X_t are independent and $\mathcal{F}_{t-1} = \sigma(X_1, \dots, X_{t-1})$, so that $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = \mathbb{E}[X_t]$.

Basic ingredient of the proof: extension of Hoeffding's lemma to conditional expectations

Lemma: X random variable s.t. $X \in [a, b]$ a.s.

Then, for all σ -algebras \mathcal{G}_j , for all $s \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{s(X - \mathbb{E}[X | \mathcal{G}_j])} | \mathcal{G}_j] = \ln (\mathbb{E}[e^{sx} | \mathcal{G}_j]) - s \mathbb{E}[x | \mathcal{G}_j] \leq \frac{s^2}{8}(b-a)^2$$

(we will discuss the proof later on... let's first prove the theorem based on this lemma.)

Proof (of the theorem):

Markov-Chebyshev bounding (= Markov's inequality after taking exponents):

$$\text{We denote } S_T = \sum_{t=1}^T X_t - \underbrace{\mathbb{E}[X_t | \mathcal{F}_{t-1}]}_{\text{martingale increments or martingale differences}}$$

(martingale = sum of
martingale increments or martingale differences)

The « probabilistic version » is about upper bounding $\mathbb{P}\{S_T > \varepsilon\}$:

$$\mathbb{P}\{S_T > \varepsilon\} = \mathbb{P}\{e^{\lambda S_T} > e^{\lambda \varepsilon}\} \leq e^{-s\varepsilon} \mathbb{E}[e^{sS_T}]$$

↑
Hölders inequality

We show by induction that $\mathbb{E}[e^{\lambda S_T}] \leq \exp\left(\frac{\lambda^2}{8} \sum_{t=1}^T (b_t - a_t)^2\right)$

- For $T=1$, true by the conditional version of Hoeffding's lemma and the fact that $S_1 = X_1 - \mathbb{E}[X_1 | \mathcal{F}_0]$ with $X_1 \in [a_1, b_1]$
- For $T-1 \rightarrow T$, where $T \geq 2$:
+ taking expectations by tower rule: $\mathbb{E} = \mathbb{E}[\mathbb{E}[\cdot | \mathcal{F}_0]]$

The extension of Hoeffding's lemma ensures that

$$\mathbb{E}[e^{\lambda(X_T - \mathbb{E}[X_T | \mathcal{F}_{T-1}])} | \mathcal{F}_{T-1}] \leq e^{\lambda^2(b_T - a_T)^2/8}$$

$$\begin{aligned} \text{so that } \mathbb{E}[e^{\lambda S_T}] &= \mathbb{E}[\mathbb{E}[e^{\lambda S_T} | \mathcal{F}_{T-1}]] \\ &= \mathbb{E}[e^{\lambda S_{T-1}} \mathbb{E}[e^{\lambda(X_T - \mathbb{E}[X_T | \mathcal{F}_{T-1}])} | \mathcal{F}_{T-1}]] \\ &\stackrel{\substack{\text{by the} \\ \text{induction} \\ \text{hypothesis}}}{\leq} e^{\lambda^2(b_T - a_T)^2/8} \times \mathbb{E}[e^{\lambda S_{T-1}}] \\ &\leq \exp\left(\lambda^2 \sum_{t \leq T} (b_t - a_t)^2/8\right) \end{aligned}$$

Substituting above: $\mathbb{P}\{S_T > \varepsilon\} \leq \inf_{\lambda > 0} \exp\left(-s\varepsilon + \lambda^2 \sum_{t \leq T} (b_t - a_t)^2/8\right)$

strictly convex function to
minimize in the exponent:
minimum achieved at λ^*
such that $\lambda^* \sum_{t \leq T} (b_t - a_t)^2/4 = \varepsilon$ (gradient vanishes)
i.e. $\lambda^* = 4\varepsilon / \sum_{t \leq T} (b_t - a_t)^2$

→ It only remains to prove the extension of Hoeffding's lemma to conditional expectations.

But first (reminder) \downarrow unconditional version

Lemma (Hoeffding) : X random variable s.t. $X \in [a, b]$ a.s.

Then $\forall s \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{s(X-\mathbb{E}X)}] = \ln \mathbb{E}[e^{sx}] - s\mathbb{E}x \leq \frac{s^2(b-a)^2}{8}$$

Proof (most elegant one I know of) :

$$\Psi(s) = \ln \mathbb{E}[e^{sx}] \text{ defined for all } s \in \mathbb{R}$$

Ψ is differentiable at each $s \in \mathbb{R}$: cf. X bounded, thus

$\eta \mapsto X e^{\eta x}$ locally dominated around s by an integrable r.v. independent of η thus $\eta \mapsto \mathbb{E}[e^{\eta x}]$ differentiable at s with derivative $\mathbb{E}[X e^{sx}]$

with

$$\Psi'(s) = \frac{\mathbb{E}[X e^{sx}]}{\mathbb{E}[e^{sx}]}$$

Similarly, Ψ is twice differentiable at each $s \in \mathbb{R}$, with:

$$\Psi''(s) = \frac{\mathbb{E}[X^2 e^{sx}] \mathbb{E}[e^{sx}] - (\mathbb{E}[X e^{sx}])^2}{(\mathbb{E}[e^{sx}])^2} = \text{Var}_{\mathbb{Q}}(X)$$

under the probability \mathbb{Q} defined by

$$\frac{d\mathbb{Q}}{dP}(w) = \frac{e^{sw}}{\mathbb{E}[e^{sx}]}$$

$$\begin{aligned} X \in [a, b] : \quad \text{Var}_{\mathbb{Q}}(X) &= \inf_{\mu \in \mathbb{R}} \mathbb{E}_{\mathbb{Q}}[(X-\mu)^2] \\ &\leq \mathbb{E}_{\mathbb{Q}}[(X - \frac{a+b}{2})^2] = \frac{(b-a)^2}{4} \end{aligned}$$

$$\text{Taylor: } \exists x \text{ s.t. } \Psi(s) = \underbrace{\Psi(0)}_{=0} + \underbrace{s\Psi'(0)}_{=0} + \frac{s^2}{2} \underbrace{\Psi''(x)}_{\leq (b-a)^2/4}$$

$$\text{i.e., } \ln \mathbb{E}[e^{sx}] \leq \frac{s^2(b-a)^2}{8}$$

Back to Hoeffding's lemma with conditional expectations:

Proof 1? Can we take the proof of Hoeffding's lemma we just saw and replace all E by $E[g]$?

$$\Psi(s) = \ln E[e^{sx}|g] \rightarrow \text{The theorem of differentiation under } E[\cdot] \text{ only requires dominated convergence, which holds true for } E[g|g].$$

as well. Thus, we also have a theorem of differentiation under $E[g|g]$:

a.s., $\Psi''(s)$ exists and equals $\Psi''(s) = \frac{E[x^2 e^{sx}|g] E[e^{sx}|g] - (E[x e^{sx}|g])^2}{(E[e^{sx}|g])^2}$

= some conditional variance under a different probability measure?

Yes, using the notion of « regular conditional probability », which always exists in our case, we could perform a change of measure again and identify a conditional variance.

But PLEASE! this is extremely heavy math... I can't inflict that to you....

Proof 2

Too bad for elegance, let's get back to the original proof of Hoeffding's (unconditional) lemma, which only relies on calculus:

$$y = x - E[x|g] \in [A, B] \quad \text{where } A = a - E[x|g] \leq 0 \quad B = b - E[x|g] \geq 0$$

one both g -measurable
and $B-A = b-a > 0$

$$y = \frac{B-y}{B-A} A + \frac{y-A}{B-A} B$$

↑ convex weights

Since $y \mapsto e^{sy}$ is convex:

$$e^{sy} \leq \frac{B-y}{B-A} e^{sA} + \frac{y-A}{B-A} e^{sB}$$

Taking $E[g]$: using $E[y|g] = 0$ and A, B g -measurable:

$$E[e^{sy}|g] \leq \frac{B}{B-A} e^{sA} - \frac{A}{B-A} e^{sB}$$

note that $s/B-A$ and $-A/B-A$ are convex weights

Now, by a function study (the very same as the one we performed in the proof of the unconditional version of Hoeffding's lemma) — or even by the latter lemma itself:

$$\forall u, v \in \mathbb{R}, \forall p \in [0, 1], \forall s \in \mathbb{R}, \ln(p e^{su} + (1-p)e^{sv}) \leftarrow \ln \text{ of expected value of } e^{sz} \text{ where } z = \begin{cases} u & \text{w.p. } p \\ v & \text{w.p. } 1-p \end{cases}$$

$$\leq s(pu + (1-p)v) + \frac{s^2}{8}(v-u)^2$$

↑ expected value of z
range is $[u, v]$

In particular,

$$\frac{B}{B-A} e^{sA} - \frac{A}{B-A} e^{sB} \leq \exp\left(s\left(\frac{BA}{B-A} - \frac{AB}{B-A}\right) + \frac{s^2}{8}(B-A)^2\right)$$

$$= \exp\left(\frac{s^2}{8}(b-a)^2\right) \quad \text{recall that a.s., } B-A = b-a$$

Summarizing:

$$\mathbb{E}[e^{sx} | \mathcal{G}] \leq \exp\left(\frac{s^2}{8}(b-a)^2\right)$$

$$= \mathbb{E}[e^{sx}] \times \exp(-s \mathbb{E}[x | \mathcal{G}])$$

Proof 3

My preferred (not only because I found by myself):

Hoeffding's lemma in its unconditional version ENTAILS the

Conditional version! This is because Hoeffding's lemma holds for all probability distributions — we should play with this fact.

For all $A \in \mathcal{G}$
s.t. $\mathbb{P}(A) > 0$, let $\mathbb{P}_A = \mathbb{P}(\cdot | A)$, the conditional distribution given the event A .

The unconditional version of Hoeffding's lemma ensures that

$$\forall \epsilon \in \mathcal{G} \text{ s.t. } \mathbb{P}(A) > 0, \quad \forall s \in \mathbb{R}, \quad \ln \mathbb{E}_A[e^{sx}] \leq s \mathbb{E}_A[x] + \frac{s^2}{8}(b-a)^2$$

Why do we consider the \mathbb{E}_A ? Because $\mathbb{E}[x | \mathcal{G}]$ is the unique \mathcal{G} -measurable random variable such that
 $\forall \epsilon \in \mathcal{G}$, $\mathbb{E}[x \mathbf{1}_A] = \mathbb{E}[\mathbb{E}[x | \mathcal{G}] \mathbf{1}_A]$
or, equivalently, $\forall \epsilon \in \mathcal{G}$ s.t. $\mathbb{P}(A) > 0$, $\mathbb{E}_A[x] = \mathbb{E}_A[\mathbb{E}[x | \mathcal{G}]]$.

Now, consider the event $H = \{ \mathbb{E}[e^{sx} | \mathcal{G}] - e^{s \mathbb{E}[x | \mathcal{G}]} e^{\frac{s^2(b-a)^2}{8}} > 0 \} \in \mathcal{G}$

We want to prove that $P(H) = 0$. We proceed by contradiction:

If we had $P(H) > 0$, then P_H would be defined and

$$\begin{aligned} E_H[e^{sx}] &= E_H[E[e^{sx}|G]] > \underset{\substack{\text{by definition of } H \\ \text{and because a function } > 0 \\ \text{on a set with probability } > 0 \\ \text{has an expectation } > 0}}{E_H[e^{sE[X|G]}]} e^{s^2(b-a)^2/8} \\ &\stackrel{\substack{> \\ (\text{unconditional}) \text{ Jensen's} \\ \text{inequality}}}{=} e^{sE_H[X]} e^{s^2(b-a)^2/8} \end{aligned}$$

which would be in contradiction with the unconditional Hoeffding's lemma:

$$\ln E_H[e^{sx}] \leq sE_H[X] + \frac{s^2}{8}(b-a)^2.$$

We conclude that $P(H) = 0$.

A final remark:

A better version of Hoeffding's lemma in its conditional form is the following one:

Hoeffding's lemma: X random variable s.t. there exists G G_j -measurable and integrable, as well as $a, b \in \mathbb{R}$ with: $G+a \leq X \leq G+b$ a.s.
Then $\forall s \in \mathbb{R}$, $\ln E[e^{sx}|G] - sE[X|G] \leq \frac{s^2}{8}(b-a)^2$

⇒ In my first statement, I only considered $G = \emptyset$.

Why can we consider general G 's? See next page.

Then we get the following extension of the Hoeffding-Azuma inequality:

i: Let (\mathcal{F}_t) be a filtration and $(X_t)_{t \geq 0}$ be a sequence of adapted r.v. such that

$\forall t, \exists G_t \mathcal{F}_{t-1}$ -measurable and integrable, $\exists a_t, b_t | a_t + G_t \leq X_t \leq b_t + G_t$

then, $\forall s \in \mathbb{R}$, w.p 1-s,

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

Why does Hoeffding's lemma hold for more general G 's?

We consider $a \leq X - G \leq b$ as

Hoeffding's lemma indicates that

$$\ln E[e^{s(X-G)} | \mathcal{G}_j] - s E[X-G | \mathcal{G}_j] \leq \frac{s^2}{8} (b-a)^2$$

We have $G \in L^1$ thus $X \in L^1$ and therefore,

$$s E[X-G | \mathcal{G}_j] = E[X | \mathcal{G}_j] - G \in L^1$$

Now, we also have $E[e^{s(X-G)} | \mathcal{G}_j] = e^{-sG} E[e^{sX} | \mathcal{G}_j]$

Hence the result by cancelling out the $-sG$ terms.

(I don't think we need here that $e^{sG} \in L^1$ or equivalently $e^{sX} \in L^1$;
do we?)

↳ maybe we do / you all seemed to say that we did...
anyway in most of our applications, G will be bounded as well!

What can we do when no convexity assumption holds?

↳ Non-convex aggregation via randomization

Example 1:

N-ary decisions
in a game (4-ary if we have to pick paths in a graph: $\rightarrow \leftarrow \uparrow \downarrow$)

1. Opponent picks state of the world y_t
2. Statistician picks action $j_t \in \{1, \dots, N\}$
3. Loss $l(j_t, y_t)$ or reward $-l(j_t, y_t)$ is encountered, both y_t and j_t are made public

Example 2: Prediction with expert advice (the «meta-statistical» framework)

↳ when the prediction space is not convex:

1. Opponent picks observation $y_t \in \mathcal{Y}$
2. Simultaneously, experts provide forecasts $f_{jt} \in \mathcal{Y}, j \in \{1, \dots, N\}$
and statistician picks forecast $\hat{j}_t \in \mathcal{Y}$
3. y_t and \hat{j}_t are revealed, losses $l(\hat{j}_t, y_t)$ and $l(f_{j_t}, y_t)$ are suffered

No convexity: \mathcal{Y} not convex [or \mathcal{Y} convex but $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
eg, $\mathcal{Y} = \{1, \dots, M\}$ in Many classification not convex in its first argument]

↳ \hat{j}_t cannot be any convex/linear prediction of the f_{jt} we wish.

Solution:

(at least,
an easy solution,
these might be
others)

Draw $J_t \in \{1, \dots, N\}$ at random

and pick $\begin{cases} \text{action } J_t \text{ (in Example 1)} \\ \text{forecast } \hat{j}_t = f_{J_t, t} \text{ (in Example 2)} \end{cases}$



General setting:

Simultaneously $\begin{cases} 1. \text{Opponent picks } \ell = (\ell_{1,t}, \dots, \ell_{N,t}) \in \mathbb{R}^N \\ 2. \text{Statistician draws } J_t \in \{1, \dots, N\} \end{cases}$

3. J_t and $(\ell_{1,t}, \dots, \ell_{N,t})$ are revealed

Aim: Minimize the regret

$$\sum_{t=1}^T l_{J_t, t} - \min_{k=1, \dots, N} \sum_{t=1}^T \ell_{k,t}$$

! The losses l_{jt} may depend on the past, i.e., on J_1, \dots, J_{t-1}

Methodology: We denote by $p_t = (p_{1t}, \dots, p_{Nt}) \in \mathcal{X}$ the probability distribution used to draw J_t , conditionally to the past

$$\text{Regret: } R_T = \sum_{t=1}^T l_{J_t, t} - \min_k \sum_{t=1}^T l_{k, t} = \sum_{t=1}^T l_{J_t, t} - \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt}$$

$$+ \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} - \min_k \sum_{t=1}^T l_{kt}$$

This can be controlled independently of the probability distributions chosen

We already learned how to control this term!

The information available at the beginning of round t is $(l_s, p_s, J_s)_{s \leq t-1}$

We denote $\mathcal{F}_{t-1} = \sigma\{(l_s, p_s, J_s)_{s \leq t-1}\}$: l_t and p_t are \mathcal{F}_{t-1} -measurable while J_t is drawn at random using an auxiliary randomization $U_t \sim U_{[q_1]}$, independent from \mathcal{F}_{t-1} .

Then: $E[l_{J_t, t} | \mathcal{F}_{t-1}] = \sum_{j=1}^N p_{jt} l_{jt}$ (J_t is not fixed by the conditioning, only its distribution p_t is)

↳ Expected regret (conditionally expected regret)

$$\bar{R}_T = \sum_{t=1}^T p_{jt} l_{jt} - \min_k \sum_{t=1}^T l_{kt}$$

We already saw that we could ensure $\bar{R}_T \leq O((M-m)\sqrt{T \ln N})$ if $l_{jt} \in [m, M]$

↳ Martingale

$$S_T = \sum_{t=1}^T l_{J_t, t} - \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt}$$

The Hoeffding-Azuma inequality

ensures that

if $l_{jt} \in [m, M] \forall j \neq t$, then, no matter which p_t were selected

$$\begin{aligned} \text{with} \\ l_{jt} &= E_{J_t} l_{jt} \\ l_{tt} &= M \text{ and} \\ l_{tt} &= m \\ l_{jt} &= \mathbb{E}_{J_t} l_{jt} \\ G_t &= \sum_j l_{jt} \end{aligned}$$

$$\mathbb{P}\{ S_T \leq (M-m)\sqrt{\frac{T}{2} \ln \frac{1}{\delta}} \} \geq 1-\delta$$

Conclusion: \bar{R}_T , with probability at least $1-\delta$,

$$R_T \leq \bar{R}_T + (M-m)\sqrt{\frac{T}{2} \ln \frac{1}{\delta}}$$

E.g. with the fully adaptive version of EWA:

$$T, \forall \epsilon > 0, \text{ with probability at least } 1-\delta, \quad R_T \leq (M-m) \sqrt{T} \left(\sqrt{\ln N} + \sqrt{\frac{1}{2} \ln \frac{1}{\delta}} \right) + (M-m)(2 + 4\sqrt{3} \ln N)$$

This is called a high probability bound; it is non-asymptotic \rightarrow Exercise: Can you get a high probability bound of the form $\forall \epsilon > 0, \text{ with probability at least } 1-\delta, \forall T, R_T \leq \dots$?

Consequence: Asymptotic almost-sure bound.

The Borel-Cantelli lemma, using $S_T = 1/T^2$,

$$\mathbb{P} \left(\limsup_{T \rightarrow \infty} \{ R_T > (M-m) \sqrt{T} \left(\sqrt{\ln N} + \sqrt{\ln T} \right) \} \right) = 0$$

limsup of events

We denote $\rho(T)$
This quantity:

$$\rho(T) \sim (M-m) \sqrt{T \ln T}$$

That is, almost-surely

$$R_T / \rho(T) > 1 \text{ for finitely many } T$$

thus $\limsup_{T \rightarrow \infty} \frac{R_T}{\rho(T)} \leq 1 \text{ a.s.}$ or equivalently,

↑
limsup of
a sequence of numbers

$$\limsup_{T \rightarrow \infty} \frac{R_T}{(M-m) \sqrt{T \ln T}} \leq 1 \text{ a.s.}$$

Exercise: [To be stated in a more detailed way on the next page]

Show that we actually have

$$\limsup_{T \rightarrow \infty} \frac{R_T}{(M-m) \sqrt{T \ln(\ln T)}} \leq C \text{ a.s.}$$

where C
is a constant

(a rate which should
be reminiscent of the law of the iterated logarithm -)

and I should have
started with that...

Note: Of course, since $E[S_T] = 0$, we have $E[R_T] = E[\bar{R}_T]$

Because we have deterministic bounds on \bar{R}_T , we get bounds on

$E[R_T]$. But this doesn't tell us much on R_T , this is

why we prefer our high-probability bounds -

Exercise

[Full Statement]

- (1) Remind yourself of Doob's martingale inequality
 (actually: inequalities - there are two of them, but we'll need only the most famous one).

- (2) Show the following MAXIMAL version of the Hoeffding-Azuma inequality:

$\forall \delta \in (0,1)$, with probability at least $1-\delta$,

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

- (3) Show that for any algorithm with expected regret \bar{R}_T less than something of order $(M-m)\sqrt{T \ln N}$, the corresponding randomized algorithm has a regret R_T such that

For all strategies of the opponent picking losses $a_t \in [m, M]$,

$$\limsup_{T \rightarrow \infty} \frac{R_T}{(M-m)\sqrt{T \ln(\ln T)}} \leq C \quad \text{a.s.}$$

where C is a universal constant (propose a numerical value).

- (4) Is this C optimal? (Consider the law of the iterated logarithm as a basis for your discussion.)

Hint for (3):

Consider the regimes $\{2^{r+1}, \dots, 2^{r+1}\}$ for $r=1, 2, \dots$ and pick $S_r = 1/r^2$ for the application of the Borel-Cantelli lemma. (cf. doubling trick!)

Note: If you can solve (3) and have interesting things to say for (4), please send me your work.

Stochastic bandits.Finitely many arms.Setting: K arms indexed by $1, 2, \dots, K$ With each arm j is associated a probability distribution π_j
(over \mathbb{R})
with an expectationAt each round $t = 1, 2, \dots$

- The decision-maker picks $I_t \in \{1, \dots, K\}$, possibly at random
- She gets a reward y_t drawn at random according to π_{I_t} (given I_t)
- This is the only feedback she gets / the only observation she has access to.

Aim:We denote by $\mu_i = E(y_i)$ the expectation of y_i (note: operator E vs. expectation E of an expression involving random variables.)Pseudo-regret $R_T = T\mu^* - E\left[\sum_{t=1}^T y_t\right]$ to be controlledwhere $\mu^* = \max_{j \in K} \mu_j$ Useful notation: $\Delta_a = \mu^* - \mu_a$ gap of arm a $\Delta_a = 0$: a is an optimal arm (there can be several of them) $\Delta_a > 0$: a is a suboptimal arm $N_a(T) = \sum_{t=1}^T \mathbf{1}_{I_t=a}$ total number of times that a is pulled.Note: * Pseudo regret R_T is a very "expected" notion of regret

$$R_T \leq \underbrace{\text{probably}}_{E\left[\max_{a=1, \dots, K} \frac{1}{t} \sum_{t=1}^t y_a - \frac{1}{t} \sum_{t=1}^t y_t\right]}$$

* Can be rewritten (later) as $R_T = \sum_{a=1}^K \Delta_a E[N_a(T)]$

Upper confidence bound [UCB] algorithm: very popular!

For $t = 1, 2, \dots, K$

- Pull arm $I_t = t$, get a reward y_t

For $t = K+1, K+2, \dots$

- Pull an arm $I_t \in \arg\max_{j \in \{1, \dots, K\}} \left\{ \hat{\mu}_{j,t-1} + \sqrt{\frac{2 \ln t}{N_j(t-1)}} \right\}$

(sel.-braking rule)
pick the element with
smallest index

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

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

always ≥ 1 since
each arm was tried
sequentially during rounds
 $1, 2, \dots, K$

- Get a reward y_t

Theorem: If the distributions y_j have supports all included in $[a_1]$, then the pseudo-regret of UCB is smaller than

$$\bar{R}_T \leq \sum_{i: \Delta_i > 0} \left(\frac{8 \ln T}{\Delta_i} + 2 \right)$$

This regret bound is obtained via the following proposition:

Proposition: If the distributions y_j have supports all included in $[a_1]$,

then

$$\forall i \text{ s.t. } \Delta_i > 0, \quad E[N_i(T)] \leq \frac{8 \ln T}{\Delta_i^2} + 2.$$

Exercise

The bounds above are called distribution-dependent because they depend heavily on the distributions y_i at hand (via the gaps $\Delta_i = \mu^* - \mu_i$).

Show the following distribution-free bound (that only

depends on the support $[q_1]$, not on the specific distributions π_i^* at hand) : for the UCB algorithm,

$$\sup_{\substack{\pi_1, \dots, \pi_K \text{ with} \\ \text{supports in } [q_1]}} \bar{R}_T \leq O(\sqrt{T \ln T}).$$

Hint: For small values of Δ_i , the bound of the Proposition can be worse than the trivial T bound...

Proof [of the theorem based on the Proposition] :

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

where by definition of the bandit model, \sim Given I_t , y_t is drawn at random according to π_{I_t}

$$E[y_t | I_t] = \mu_{I_t}$$

thus (by the tower rule)

$$\begin{aligned} E[y_t] &= E[\mu_{I_t}] \\ &= \sum_j \Delta_j E[1_{\{I_t=j\}}] \end{aligned}$$

Summing over t :

$$E\left[\sum_{t=1}^T y_t\right] = \sum_{j=1}^K \Delta_j E[N_j(T)]$$

and (in view of $T = \sum_j E[N_j(T)]$)

$$\begin{aligned} \bar{R}_T &= \sum_j (\mu^* - \mu_j) E[N_j(T)] = \sum_{j=1}^K \Delta_j E[N_j(T)] \\ &= \sum_{j: \Delta_j > 0} \Delta_j E[N_j(T)] \end{aligned} \quad \begin{array}{l} \text{it suffices} \\ \text{to consider} \\ \text{the suboptimal} \\ \text{arms...} \end{array}$$

We conclude by substituting $E[N_j(T)] \leq \frac{8 \ln T}{\Delta_j^2} + 2$ and by bounding $2\Delta_j \leq 2$.

Note: Keep in mind the rewriting as we will often use it!

$$\begin{aligned} \bar{R}_T &= T\mu^* - E\left[\sum_{t=1}^T y_t\right] \\ &= \sum_{a=1}^K \Delta_a E[N_a(T)] \end{aligned}$$

Proof [of the Proposition]: We fix an optimal arm $a^* \in \{1, \dots, K\}$, i.e. s.t. $\mu_{a^*} = \mu^{*+}$

→ It will show why this algorithm is called UCB:

Because $\hat{\mu}_{j,t-1} + \sqrt{\frac{2 \ln t}{N_j(t-1)}}$ will indeed appear as an upper confidence bound on μ_j

estimate based on the raw performance
↳ exploration of the results

larger for arms j not much sampled so far
↳ forces some exploration

The UCB algorithm realizes some compromise / trade off between exploitation & exploration

Later on we compare these state wants to life Hoeffding-Azuma inequality

[G]

S

LEMMA:

$\forall j, \forall t \geq j$ (so that $N_j(t) \geq 1$)

$\hat{\mu}_j = \frac{\sum_{i=1}^t \mathbb{I}_{I_i=j}}{N_j(t)}$ is supported by $\{1, \dots, t\}$ and $\hat{\mu}_j > \mu_j$ if $\mu_j < \hat{\mu}_j$

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

$$\mathbb{P}\left\{\mu_j > \hat{\mu}_{j,t} - \sqrt{\frac{\ln(1/\delta)}{2N_j(t)}}\right\} \geq 1 - \delta.$$

Or

By symmetry: $\forall \delta \in (0, 1)$,

$$\mathbb{P}\left\{\mu_j < \hat{\mu}_{j,t} + \sqrt{\frac{\ln(1/\delta)}{2N_j(t)}}\right\} \geq 1 - \delta$$

→ Application :

$$N_i(T) = 1 + \sum_{t=K+1}^T \mathbb{I}_{I_t=i}$$

We show below that $t \geq K+1$ and $I_t=i$ entails one of the following:

(i) $\hat{\mu}_{i,t-1} > \mu_i + \sqrt{\frac{2 \ln t}{N_i(t-1)}}$ [$\mu_i <$ lower confidence bound]

(ii) $\hat{\mu}_{a^*,t-1} < \mu^{*+} - \sqrt{\frac{2 \ln t}{N_{a^*}(t-1)}}$ [$\mu^{*+} >$ upper confidence bound]

(iii) $N_i(t-1) \leq \frac{8 \ln t}{\Delta^2}$ [i not played often yet]

Indeed, we would otherwise have

$$\hat{\mu}_{\alpha^*, t-1} + \sqrt{\frac{2 \ln t}{N_{\alpha^*}(t-1)}} \geq \mu^*$$

$= \mu_i + \Delta$

negation of (ii)

$$> \mu_i + 2 \sqrt{\frac{2 \ln t}{N_i(t-1)}} \quad \left. \begin{array}{l} \text{definition of } \Delta_i \\ \text{the negation of (iii)} \\ \text{is } \Delta_i^2 > 8 \ln t / N_i(t-1) \end{array} \right\}$$

$$\geq \hat{\mu}_{i, t-1} + \sqrt{\frac{2 \ln t}{N_i(t-1)}} \quad \text{negation of (i)}$$

\Rightarrow inequality between these quantities would contradict $I_t = i$, that is,
 $\in \arg \max_j \{ \hat{\mu}_j + \sqrt{2 \ln t / N_j(t-1)} \}$

Thus, $E[N_i(\tau)] \leq 1 + \sum_{t=K+1}^T \mathbb{P}(\hat{\mu}_{i, t-1} > \mu_i + \sqrt{\frac{2 \ln t}{N_i(t-1)}}) + \mathbb{P}(\hat{\mu}_{\alpha^*, t-1} < \mu^* - \sqrt{\frac{2 \ln t}{N_{\alpha^*}(t-1)}})$

each $\leq t \delta$
where $\delta = \Delta_i^2 / t^4$

\downarrow

$+ E \left[\sum_{t=K+1}^T \mathbb{1}_{\{I_t = i \text{ & } N_i(t-1) \leq 8 \ln t / \Delta_i^2\}} \right] \quad \downarrow 8 \ln t \leq 8 \ln T$

$\leq 1 + 2 \sum_{t=K+1}^T t^{-3} + E \left[\sum_{t=K+1}^T \mathbb{1}_{\{N_i(t-1) \leq 8 \ln T / \Delta_i^2 \text{ & } I_t = i\}} \right]$

$\leq 1 + 2 \sum_{t=K+1}^T t^{-3} + \underbrace{\left(\frac{8 \ln T}{\Delta_i^2} + 1 \right)}_{\text{deterministically upper bounded by}} - 1$

\downarrow

$\text{as } I_t = i \text{ only if } N_i(t-1) \leq \frac{8 \ln T}{\Delta_i^2} + 1$

$\text{thus only if } N_i(t) \leq \frac{8 \ln T}{\Delta_i^2} + 1$

-1

$\text{because } I_t = i$

$\text{is not included in the sum}$

\downarrow

$\int_1^{+\infty} t^{-3} dt = [-t^{-2}]_1^{+\infty} = 1$

$\text{so that the total sum is controlled by this number}$

$\sum_{t=1}^T \mathbb{1}_{\{I_t = i\}} = N_i(t)$

$\sum_{t=K+1}^T \mathbb{1}_{\{I_t = i\}}$

thus:

$$E[N_i(\tau)] \leq \frac{8 \ln T}{\Delta_i^2} + 2$$

Proof of the lemma (Hoeffding-Azuma inequality with a random number of summands):

Let

$$Z_t = \sum_{s=1}^t (Y_s - \mu_a) \mathbf{1}_{\{I_s=a\}}; \quad \text{we successively prove:}$$

(0) $(Z_t)_{t \geq 0}$ is a martingale wrt. $(\mathcal{F}_t)_{t \geq 0} = (\sigma(Y_1, \dots, Y_t))$

where $\mathcal{F}_0 = \{\emptyset, \Omega\}$ trivial σ -algebra

Indeed: each I_t is \mathcal{F}_{t-1} -measurable (picked based only on past payoffs)

thus Z_t is \mathcal{F}_t -adapted

Showing that it is a martingale amounts to showing

$$\mathbb{E}[(Y_t - \mu_a) \mathbf{1}_{\{I_t=a\}} \mid Y_1, \dots, Y_{t-1}] = 0 \quad \text{a.s.}$$

but since I_t is \mathcal{F}_{t-1} -measurable, this quantity equals

$$\mathbb{E}[(Y_t - \mu_a) \mathbf{1}_{\{I_t=a\}} \mid Y_1, \dots, Y_{t-1}, I_t]$$

$$= (\mathbb{E}[Y_t \mid I_t, Y_1, \dots, Y_{t-1}] - \mu_a) \mathbf{1}_{\{I_t=a\}}$$

$$= (\mu_{I_t} - \mu_a) \mathbf{1}_{\{I_t=a\}} = 0 \quad \text{a.s., as desired}$$

Given I_t , thus by the
bandit model, its conditional
expectation equals μ_{I_t}

Then: (try to prove these statements by yourself, as an exercise for the next lesson):

(1) For all $x \in \mathbb{R}$, $(M_t) = \left(\exp(x Z_t - \frac{x^2}{8} N_a(t)) \right)_{t \geq 0}$ is an $(\mathcal{F}_t)_{t \geq 0}$ adapted supermartingale

↳ in particular $\mathbb{E}[M_t] \leq 1$ for all t

(2) $\forall \varepsilon > 0, \forall l \geq 1, \quad \mathbb{P}\{Z_t \geq \varepsilon \text{ and } N_a(t) = l\} \leq e^{-\frac{2\varepsilon^2}{l}}$

(3) From these we will conclude.