# Part 1: The Hoeffding–Azuma inequality

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 $\mathcal{F}_t$–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 ( « probabilistic version »)

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

or ( « statistical version », totally equivalent)

$$\forall \delta \in (0,1), \qquad \text{with probability at least } 1-\delta, \qquad \sum_{t=1}^{T} X_t - \sum_{t=1}^{T} \mathbb{E}[X_t \mid \mathcal{F}_{t-1}] \leq \sqrt{\frac{\sum_{t=1}^{T}(b_t - a_t)^2}{2} \ln\frac{1}{\delta}}$$

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 $\mathbb{E}[X_t \mid \mathcal{F}_{t-1}] = \mathbb{E}[X_t]$.

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 $\sigma$-algebras $\mathcal{G}$, for all $s \in \mathbb{R}$,

$$\ln \mathbb{E}\left[ e^{s(X - \mathbb{E}[X \mid \mathcal{G}])} \mid \mathcal{G} \right] = \ln\left( \mathbb{E}[e^{sX} \mid \mathcal{G}] \right) - s\,\mathbb{E}[X \mid \mathcal{G}] \leq \frac{s^2}{8}(b-a)^2$$

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

Proof (of the theorem):
Markov–Chernov bounding (= Markov's inequality after taking exponents):

We denote $S_T = \sum_{t=1}^{T} \underbrace{X_t - \mathbb{E}[X_t \mid \mathcal{F}_{t-1}]}_{}$

( $\underset{\uparrow}{\text{martingale}}$ = sum of $\underset{\uparrow}{\text{martingale increments}}$ )
or martingale differences

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

$$\mathbb{P}\{S_T > \varepsilon\} = \mathbb{P}\{e^{sS_T} > e^{s\varepsilon}\} \underset{\text{Markov's inequality}}{\leq} e^{-s\varepsilon}\,\mathbb{E}[e^{sS_T}]$$

$$\uparrow \;\; \forall s > 0$$

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

- 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 $X_1 \in [a_1, b_1]$

           + taking expectations by tower rule: $\mathbb{E} = \mathbb{E}[\mathbb{E}[\,\cdot\,|\mathcal{F}_0]]$

- For $T-1 \to T$, where $T \geq 2$:

     The extension of Hoeffding's lemma ensures that

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

so that $\mathbb{E}[e^{sS_T}] = \mathbb{E}\left[\mathbb{E}[e^{sS_T}|\mathcal{F}_{T-1}]\right]$

$$= \mathbb{E}\left[e^{sS_{T-1}}\,\mathbb{E}[e^{s(X_T - \mathbb{E}[X_T|\mathcal{F}_{T-1}])}|\mathcal{F}_{T-1}]\right]$$

$$\leq e^{s^2(b_T - a_T)^2/8} \times \mathbb{E}[e^{sS_{T-1}}]$$

by the induction hypothesis $\Bigg($

$$\leq \exp\left(s^2 \sum_{t \leq T}(b_t - a_t)^2/8\right)$$

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

strictly convex function to minimize in the exponent: minimum achieved at $s^*$

$$= \exp\left(-2\varepsilon^2 \Big/ \sum_{t \leq T}(b_t - a_t)^2\right).$$

         Such that $s^*\sum_{t \leq T}(b_t - a_t)^2/4 = \varepsilon$   (gradient vanishes)

         ie $s^* = 4\varepsilon \Big/ \sum_{t \leq T}(b_t - a_t)^2$

$\longrightarrow$ It only remains to prove the extension of Hoeffding's lemma to conditional expectations.

But first (reminder) ⌐ unconditional version

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

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

$$\ln \mathbb{E}\left[e^{s(X-\mathbb{E}X)}\right] = \ln \mathbb{E}\left[e^{sX}\right] - s\,\mathbb{E}[X] \le \frac{s^2}{8}(b-a)^2$$

<u>Proof</u>   (most elegant one I know of) :

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

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

$$\Psi''(s) = \frac{\mathbb{E}[X^2 e^{sX}]\,\mathbb{E}[e^{sX}] - \left(\mathbb{E}[X e^{sX}]\right)^2}{\left(\mathbb{E}[e^{sX}]\right)^2} = \mathrm{Var}_Q(X)$$

under the probability $Q$ defined by

$$\frac{dQ}{d\mathbb{P}}(\omega) = \frac{e^{sX(\omega)}}{\mathbb{E}[e^{sX}]}$$

$X \in [a,b]$ :        $\mathrm{Var}_Q(X) = \inf_{\mu \in \mathbb{R}} \mathbb{E}_Q\left[(X-\mu)^2\right]$

$$\le \mathbb{E}_Q\left[\left(X - \tfrac{a+b}{2}\right)^2\right] \le \frac{(b-a)^2}{4}$$

Taylor:   $\exists x$ s.t.   $\Psi(s) = \underbrace{\Psi(0)}_{=0} + s\,\underbrace{\Psi'(0)}_{=\mathbb{E}[X]} + \frac{s^2}{2}\underbrace{\Psi''(x)}_{\le (b-a)^2/4}$

Cf. $\Psi$ is actually C² smooth

i.e,

$$\ln \mathbb{E}\left[e^{sX}\right] \le s\,\mathbb{E}[X] + \frac{s^2}{8}(b-a)^2$$

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 $E[\cdot | \mathcal{G}]$?

$\Psi(s) = \ln E[e^{sX} | \mathcal{G}] \longrightarrow$ The theorem of differentiation under $E[\cdot]$ only requires dominated convergence, which holds true for $E[\cdot | \mathcal{G}]$ as well. Thus we also have a theorem of differentiation under $E[\cdot | \mathcal{G}]$:

a.s., $\Psi''(s)$ exists and equals

$$\Psi''(s) = \frac{E[X^2 e^{sX} | \mathcal{G}] \, E[e^{sX} | \mathcal{G}] - \left(E[X e^{sX} | \mathcal{G}]\right)^2}{\left(E[e^{sX} | \mathcal{G}]\right)^2}$$

$$= \text{some conditional variance under a different probability measure?}$$

Yes, see details in some pages.

However, there are two other proofs that I find more elementary:

Proof 2.    Too bad for elegance, let's get back to the original proof of Hoeffding's (unconditional) lemma, which only relies on calculus:

$$y = X - E[X | \mathcal{G}] \in [A, B] \qquad \text{where } \begin{array}{l} A = a - E[X | \mathcal{G}] \leq 0 \\ B = b - E[X | \mathcal{G}] \geqslant 0 \end{array}$$

are both $\mathcal{G}$-measurable and $B - A = b - a > 0$

$$y = \underbrace{\frac{B - y}{B - A}}_{} A + \underbrace{\frac{y - A}{B - A}}_{} B$$
convex weights

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

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

Taking $E[\cdot | \mathcal{G}]$:    using $E[y | \mathcal{G}] = 0$ and $A, B$ $\mathcal{G}$-measurable:

$$E[e^{sy} | \mathcal{G}] \leqslant \frac{B}{B - A} e^{sA} - \frac{A}{B - A} e^{sB}$$

← note that $B/B-A$ 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 underlined version of Hoeffding's lemma ) — or even by the latter lemma itself :

$$\forall u, v \in \mathbb{R}, \quad \forall p \in [0,1], \quad \ln\left(p e^{su} + (1-p) e^{sv}\right) \leftarrow \ln \text{ of expected value of}$$
$$\forall s \in \mathbb{R}, \qquad \leq s\left(p u + (1-p) v\right) \qquad e^{sZ} \text{ where}$$
$$+ \frac{s^2}{8}\left(v - u\right)^2 \qquad Z = \begin{cases} u & \text{w.p. } p \\ v & 1-p \end{cases}$$

expected value of $Z$

range is $[u, v]$

In particular,

$$\frac{B}{B-A} e^{sA} - \frac{A}{B-A} e^{sB} \leq \exp\left( s\left(\frac{BA}{B-A} - \frac{AB}{B-A}\right) + \frac{s^2}{8}(B-A)^2 \right)$$
$$= \exp\left( \frac{s^2}{8}(b-a)^2 \right) \qquad \swarrow \text{ recall that a.s.,} \quad B-A = b-a$$

Summarizing:
$$\mathbb{E}\left[ e^{sY} \mid \mathcal{G} \right] \leq \exp\left( \frac{s^2}{8}(b-a)^2 \right)$$
$$= \mathbb{E}\left[ e^{sX} \mid \mathcal{G} \right] \times \exp\left( -s\, \mathbb{E}[X \mid \mathcal{G}] \right)$$

**Proof 3**   My preferred   ( not only because I found it 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   $\mathbb{P}_A = \mathbb{P}(\cdot \mid A)$,     the conditional distribution given the event $A$.

The unconditional version of Hoeffding's lemma ensures that
$$\forall A \in \mathcal{G} \text{ s.t. } \mathbb{P}(A) > 0, \qquad \forall s \in \mathbb{R}, \qquad \ln \mathbb{E}_A\left[ e^{sX} \right] \leq s\, \mathbb{E}_A[X] + \frac{s^2}{8}(b-a)^2$$

Why do we consider the $\mathbb{E}_A$ ?   Because $\mathbb{E}[X \mid \mathcal{G}]$ is the unique $\mathcal{G}$-measurable random variable such that     $\forall A \in \mathcal{G}, \qquad \mathbb{E}[X \mathbb{1}_A] = \mathbb{E}\left[ \mathbb{E}[X \mid \mathcal{G}] \mathbb{1}_A \right]$

or, equivalently,     $\forall A \in \mathcal{G}$ s.t. $\mathbb{P}(A) > 0$,   $\mathbb{E}_A[X] = \mathbb{E}_A\left[ \mathbb{E}[X \mid \mathcal{G}] \right]$.

Now, consider the random variable     $H = e^{s\mathbb{E}[X \mid \mathcal{G}]} e^{s^2(b-a)^2/8} - \mathbb{E}[e^{sX} \mid \mathcal{G}]$

We want to prove that $H \geqslant 0$ a.s.

$H$ is $\mathcal{G}$-measurable $\#$ thus suffices to show that for all $A \in \mathcal{G}$ with $\mathbb{P}(A) > 0$,

$$\mathbb{E}_A[H] \geqslant 0 \qquad \text{that is,} \qquad \mathbb{E}[H \mathbb{1}_A] \geqslant 0$$

Indeed,

$$\mathbb{E}_A[H] = e^{s^2(b-a)^2/8} \, \mathbb{E}_A\!\left[e^{s \mathbb{E}[X|\mathcal{G}]}\right] - \mathbb{E}_A\!\left[\mathbb{E}[e^{sX}|\mathcal{G}]\right]$$

$$= e^{s^2(b-a)^2/8} \, \mathbb{E}_A\!\left[e^{s\,\mathbb{E}[X|\mathcal{G}]}\right] - \mathbb{E}_A\!\left[e^{sX}\right]$$

$$\underset{\text{(Jensen)}}{\geqslant} e^{s^2(b-a)^2/8}\, e^{s\,\mathbb{E}_A[X]} - \mathbb{E}_A\!\left[e^{sX}\right] \qquad \geqslant 0$$

(unconditional version of Hoeffding's lemma)

---

Proof 1     Let's get back to it. In what follows all expectations relative to the original probability distribution $\mathbb{P}$ will be denoted by $\mathbb{E}$, and expectations under alternative distributions $\mathbb{Q}$ will be denoted by $\mathbb{E}_{\mathbb{Q}}$.

(1)    Consider the random variable $\quad L_s = \dfrac{e^{sX}}{\mathbb{E}[e^{sX}|\mathcal{G}]} \qquad \geqslant 0$

Since $\quad \mathbb{E}[L_s] = \mathbb{E}\big[\mathbb{E}[L_s|\mathcal{G}]\big] = 1, \qquad L_s$ is a density

We define the probability $\mathbb{Q}_s$ as : $\quad \dfrac{d\mathbb{Q}_s}{d\mathbb{P}} = L_s$

(2)    We show that $\quad \psi'(s) = \dfrac{\mathbb{E}[X e^{sX}|\mathcal{G}]}{\mathbb{E}[e^{sX}|\mathcal{G}]} \quad$ also equals $\quad \mathbb{E}_{\mathbb{Q}_s}[X|\mathcal{G}]$

$$\psi''(s) = \left[\text{see expression some pages ago}\right] \qquad \mathrm{Var}_{\mathbb{Q}_s}(X|\mathcal{G})$$

To do so, it suffices to prove that for all bounded random variables $Z$ (we'll pick $Z = X$ and $Z = X^2$), we have :

$$\mathbb{E}[Z L_s |\mathcal{G}] = \mathbb{E}_{\mathbb{Q}_s}[Z|\mathcal{G}]$$

or equivalently, that $\forall A \in \mathcal{G}, \qquad \mathbb{E}\big[\mathbb{E}[Z L_s|\mathcal{G}]\,\mathbb{1}_A\big] = \mathbb{E}\big[\mathbb{E}_{\mathbb{Q}_s}[Z|\mathcal{G}]\,\mathbb{1}_A\big]$
(since both sides are $\mathcal{G}$-measurable)

But : $\quad \mathbb{E}\big[\mathbb{E}[Z L_s|\mathcal{G}]\,\mathbb{1}_A\big] = \mathbb{E}[Z L_s \mathbb{1}_A] = \mathbb{E}_{\mathbb{Q}_s}[Z \mathbb{1}_A]$

$\qquad\qquad\qquad\qquad\qquad\uparrow \qquad\qquad\qquad\qquad \uparrow$
one characterization    $L_s$ is $\frac{d\mathbb{Q}_s}{d\mathbb{P}}$
of $\mathbb{E}[\cdot|\mathcal{G}]$

and on the other end

$$\mathbb{E}\left[\; \mathbb{E}_{Q_S}[Z|\mathcal{G}]\; \mathbb{1}_A\;\right]$$

$\quad\quad$ as $\quad \mathbb{E}[L_S|\mathcal{G}] = 1$ as
$\quad\quad$ by definition of $L_S$

$$=\;\mathbb{E}\left[\;\mathbb{E}_{Q_S}[Z|\mathcal{G}]\;\;\mathbb{E}[L_S|\mathcal{G}]\;\;\mathbb{1}_A\;\right]$$

$\mathbb{E}_{Q_S}[Z|\mathcal{G}]$ and $\mathbb{1}_A$
$\mathcal{G}$-measurable and
can go inside $\mathbb{E}[\;|\mathcal{G}]$

$$=\;\mathbb{E}\left[\;\mathbb{E}\left[\;L_S\;\mathbb{E}_{Q_S}[Z|\mathcal{G}]\;\mathbb{1}_A\;\Big|\;\mathcal{G}\right]\right]$$

"tower rule"

$$=\;\mathbb{E}\left[\;L_S\;\;\mathbb{E}_{Q_S}[Z|\mathcal{G}]\;\;\mathbb{1}_A\;\right]$$

$$L_S = \frac{dQ_S}{dP}$$

$$=\;\mathbb{E}_{Q_S}\left[\;\mathbb{E}_{Q_S}[Z|\mathcal{G}]\;\mathbb{1}_A\;\right]$$

$$=\;\cancel{\mathbb{E}_{Q_S}[\mathbb{E}_{Q_S}[Z|\mathcal{G}]]}\;=\;\mathbb{E}_{Q_S}[Z\,\mathbb{1}_A]$$

$\cancel{\mathbb{1}_A\;\mathcal{G}\text{-measurable}}$
by a characterization of $\mathbb{E}[\;|\mathcal{G}]$

_(left margin, vertical):_ That's the real trick, but it's not a trick, it's what is called Bayes' formula for conditional expectations.

which concludes the proof of (2).

(3) $\quad\quad \psi'(0) = \mathbb{E}[X|\mathcal{G}]\quad\quad$ (clear)

and for all $x$, $\quad\quad \psi''(x) = \mathrm{Var}_{Q_x}(X|\mathcal{G}) \;\leq\; \dfrac{(b-a)^2}{4}$

$\uparrow$ (we prove that below)

so that we may conclude as in the case of the unconditional Hoeffding's lemma.

Indeed, $\quad \forall c \in \mathbb{R},\quad \mathbb{E}_{Q_2}\left[(X-c)^2|\mathcal{G}\right] = \mathbb{E}_{Q_2}\left[\left(X - \mathbb{E}_{Q_2}[X|\mathcal{G}] + \mathbb{E}_{Q_2}[X|\mathcal{G}] - c\right)^2 \Big|\mathcal{G}\right]$

$$=\;\underbrace{\mathbb{E}_{Q_2}\left[\left(X - \mathbb{E}_{Q_2}[X|\mathcal{G}]\right)^2|\mathcal{G}\right]}_{\overset{\text{def.}}{=}\;\mathrm{Var}_{Q_x}(X|\mathcal{G})}\; +\; 2 \times 0 \;+\; \underbrace{\cdots}_{\geq 0}$$

and we take
$$c = \frac{b+a}{2}\quad \text{and use}\quad X \in [a,b]$$

to get the a.s. bound $\quad (X-c)^2 \leq \dfrac{(b-a)^2}{4}$

A **final** remark:      Dealing with non-constant but predictable ranges in the Hoeffding-Azuma inequality

     ↳ Sometimes useful to get slightly better constants

## Hoeffding's lemma:    extension #1

Setting:   $X$ random variable s.t. there exists a bounded and $G_j$-measurable random variable $G$, as well as $a, b \in \mathbb{R}$ with: $G + a \leq X \leq G + b$

Then:
$$\forall s \in \mathbb{R}, \quad \ln \mathbb{E}[e^{sX} | G_j] \leq s\, \mathbb{E}[X|G_j] + \frac{s^2}{8}(b-a)^2$$

Rk:   $X$ bounded as well, we get this statement from the first statement by considering   $a \leq X - G \leq b$

## extension #2

Setting:   What about when there exist $U, V$ two $G_j$-measurable random variables with $U \leq X \leq V$ and $X$ bounded (for $e^{sX}$ to be $\mathbb{L}^1$)?

An inspection of Proof 2 reveals that one can prove
$$\forall s \in \mathbb{R}, \quad \ln \mathbb{E}[e^{sX} | G_j] \leq s\, \mathbb{E}[X|G_j] + \frac{s^2}{8}(V-U)^2$$

Note:   To state our extension of the Hoeffding-Azuma inequality, we will need a constant bound on $V-U$:
$$V - U \leq \Delta \quad \text{where} \quad \Delta \in \mathbb{R}^+$$

But actually
$$\left\{ \begin{array}{l} U \leq X \leq V \\ V - U \leq \Delta \in \mathbb{R}^+ \end{array} \right. \quad \text{entail} \quad \frac{U+V}{2} - \frac{\Delta}{2} \leq X \leq \frac{U+V}{2} + \frac{\Delta}{2}$$

So we're back to extension #1
(in particular, $(U+V)/2$ is bounded)

## Hoeffding-Azuma inequality

Let $(\mathcal{F}_t)$ be a filtration and $(X_t)$ be a sequence of adapted random variables such that

(1)   $\forall t, \quad \exists\, G_t \quad \mathcal{F}_{t-1}$-measurable and bounded $\left.\right\}$ with $G_t + a_t \leq X_t \leq G_t + b_t$
       $\exists\, a_t, b_t \in \mathbb{R}$

possibly following from

(2)   $\forall t, \quad \exists\, U_t, V_t \quad \mathcal{F}_t$-measurable $\left.\right\}$ with $\left\{ \begin{array}{l} V_t - U_t \leq \Delta_t \\ U_t \leq X_t \leq V_t \end{array} \right.$
            bounded
       $\exists\, \Delta_t \in \mathbb{R}^+$

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

with probability at least $1-\delta$,

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

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

# Part 2: Non-convex aggregation via randomization

What can we do when no convexity assumption holds?

↳    Non - convex  aggregation  via  randomization.

Example 1:     N-ary decisions     ( 4-ary if we have to pick paths in a
in a game                  graph:  → ← ↑ ↓ )

( binary if  accept/reject actions)

1. Opponent picks state of the world $y_t$ ⎫
⎬ simultaneously
2. Statistician picks action $j_t \in \{1, ... N\}$ ⎭

3. Loss $\ell(j_t, y_t)$ or reward $-\ell(j_t, y_t)$ 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 $f_{j,t} \in \mathcal{Y}$, $j \in \{1...N\}$

and statistician picks forecast $\hat{y}_t \in \mathcal{Y}$

3. $y_t$ and $\hat{y}_t$ are revealed,    losses $\ell(\hat{y}_t, y_t)$ and $\ell(f_{j,t}, y_t)$ are suffered

No convexity:    $\mathcal{Y}$ not convex [ OR $\mathcal{Y}$ convex but $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$
eg, $\mathcal{Y} = \{1...M\}$ in N-ary     not convex in its first argument ]
classification

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

Solution:    Draw $J_t \in \{1, ... N\}$ at random
(at least,
an easy solution,    and    pick ⎰ action $J_t$ (in Example 1)
there might be         ⎱ forecast $\hat{y}_t = f_{J_t, t}$ (in Example 2)
others)

⎰

General setting:    Simultaneously ⎰ 1. Opponent picks $\ell_t = (\ell_{1t}, ... \ell_{Nt}) \in \mathbb{R}^N$
                           ⎱ 2. Statistician draws $J_t \in \{1, ... N\}$

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

Aim:    Minimize the regret    $\sum_{t=1}^{T} \ell_{J_t, t} - \min_{k=1...N} \sum_{t=1}^{T} \ell_{kt}$

◊ The losses $\ell_{j,t}$ may depend on the past, ie, on $J_1, \dots J_{t-1}$

<u>Methodology</u>: We denote by $P_t = (p_{1,t} \dots p_{N,t}) \in \mathcal{X}$ the probability distribution used to draw $J_t$, conditionally to the past

Regret : $R_T = \sum\limits_{t=1}^{T} \ell_{J_t, t} - \min\limits_{k} \sum\limits_{t=1}^{T} \ell_{k,t} = \left[ \sum\limits_{t=1}^{T} \ell_{J_t, t} - \sum\limits_{t=1}^{T} \sum\limits_{j=1}^{N} p_{j,t} \ell_{j,t} \right]$

$\qquad\qquad\qquad\qquad + \left[ \sum\limits_{t=1}^{T} \sum\limits_{j} p_{j,t} \ell_{j,t} - \min\limits_{k} \sum\limits_{t=1}^{T} \ell_{k,t} \right]$

This can be Controlled independently of the probability distributions chosen

already learned how to control this term! we denote it by $\bar{R}_T$ below

The information available at the beginning of round $t$ is $(\ell_s, p_s, J_s)_{s \leq t-1}$
We denote $\mathcal{F}_{t-1} = \sigma\{ (\ell_s, p_s, J_s)_{s \leq t-1} \}$ : $\ell_t$ and $p_t$ are $\mathcal{F}_{t-1}$ — measurable while $J_t$ is drawn at random using an auxiliary randomization $U_t \sim \mathcal{U}_{[0,1]}$, independent from $\mathcal{F}_{t-1}$.

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

↳ Expected regret (conditionally expected regret)    $\bar{R}_T = \sum\limits_{t,j} p_{j,t} \ell_{j,t} - \min\limits_{k} \sum\limits_{t=1}^{T} \ell_{k,t}$

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

↳ Martingale    $S_T = \sum\limits_{t=1}^{T} \ell_{J_t, t} - \sum\limits_{t=1}^{T} \sum\limits_{j} p_{j,t} \ell_{j,t}$

The Hoeffding_Azuma inequality    ensures that

with $X_t = \ell_{J_t,t}$
$a_t = m$ and $b_t = M$ and $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = \sum_j p_j \ell_j$

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

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

<u>Conclusion</u>: $\forall \delta$, 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:

$\forall_T, \forall \delta \in (0,1)$, with probability at least $1-\delta$, $\qquad R_T \leq (M-m)\sqrt{T}\left(\sqrt{\ln N} + \sqrt{\frac{1}{2}\ln\frac{1}{\delta}}\right)$

$$+ (M-m)\left(2 + \frac{4}{3}\ln N\right)$$

This is called a
high probability bound; it is non-asymptotic $\rightarrow$ $\underline{\text{Exercise}}$: Can you get
a high probability
bound of the form:
$\forall \delta \in (0,1)$, with proba $\geq 1-\delta$,
$\forall_T, R_T \leq \dots$?

$\underline{\text{Consequence}}$ :   Asymptotic almost-sure bound.

The Borel-Cantelli lemma, using $\delta_T = 1/T^2$,
ensures that

$$\mathbb{P}\left(\limsup \left\{ R_T > (M-m)\sqrt{T}\left(\sqrt{\ln N} + \sqrt{\ln T}\right) \atop + (M-m)(2 + 4/3 \ln N) \right\}\right) = 0$$

$\uparrow$
limsup of events

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

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

That is, almost-surely

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

thus   $\limsup\limits_{T\to+\infty} \dfrac{R_T}{\rho(T)} \leq 1$   a.s.   or equivalently,

$\uparrow$
limsup of
a sequence of numbers

$$\limsup\limits_{T\to+\infty} \frac{R_T}{(M-m)\sqrt{T\ln T}} \leq 1 \quad \text{a.s.}$$

$\underline{\text{Exercise}}$:   [To be stated in a more detailed way on the next page]

Show that we actually have   $\limsup\limits_{T\to+\infty} \dfrac{R_T}{(M-m)\sqrt{T\ln(\ln T)}} \leq C$   a.s.

where C
is a constant

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

and I should have
started with that...

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

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

$\mathbb{E}[R_T]$.   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:

$$\forall \delta \in (0,1), \quad \text{with probability at least } 1-\delta, \qquad \max_{t \leq T} \left\{ \sum_{s=1}^{t} X_s - \sum_{s=1}^{t} \mathbb{E}[X_s \mid \mathcal{F}_{s-1}] \right\}$$

$$\leq \sqrt{\frac{\sum_{t=1}^{T}(b_t-a_t)^2}{2} \ln \frac{1}{\delta}}$$

(3)     Show that for any algorithm with expected regret $\overline{R}_T$ 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
        $\ell_{j,t} \in [m,M]$,

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

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

# Part 3: Introduction to stochastic bandits

Stochastic bandits.                    Finitely many arms.

Setting:

$K$ arms indexed by $1, 2, \ldots K$

With each arm $j$ is associated a probability distribution $\nu_j$ (over $\mathbb{R}$) with an expectation.

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

- The decision-maker picks $I_t \in \{1, \ldots K\}$, possibly at random

- She gets a reward $Y_t$ drawn at random according to $\nu_{I_t}$ (given $I_t$)

- This is the only feedback she gets / the only observation she has access to.

Aim:

We denote by $\mu_i = E(\nu_i)$ the expectation of $\nu_i$

(note: operator $E$ vs. expectation $\mathbb{E}$ of an expression involving random variables. )

Pseudo-regret $\overline{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^{T} Y_t\right]$      to be controlled

where $\mu^* = \max_{j \leqslant K} \mu_j$

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

$N_a(T) = \sum_{t=1}^{T} \mathbb{1}_{\{I_t = a\}}$      total number of times that $a$ is pulled.

Note:    * Pseudo regret $\overline{R}_T$ is a very "expected" notion of regret

$\overline{R}_T \quad \leqslant_{probably} \quad \mathbb{E}\left[\max_{a=1, \ldots K} \sum_t \ldots - \sum_{t=1}^{T} Y_t\right]$

* Can be rewritten (see later) as      $\overline{R}_T = \sum_{a=1}^{K} \Delta_a \, \mathbb{E}[N_a(T)]$

Upper confidence bound [UCB] algorithm:                    Very popular!

For $t = 1, 2 \ldots K$

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

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

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

(tie - breaking rule: pick the element with smallest index)

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

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 $\geq 1$ since each arm was tried sequentially during rounds $1, 2 \ldots K$

  - Get a reward $y_t$

Theorem:      If the distributions $y_j$ have supports all included in $[0,1]$, then the pseudo-regret of UCB is smaller than

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

This regret bound is obtained via the following proposition:

Proposition:      If the distributions $y_j$ have supports all included in $[0,1]$, then

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

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 $[0,1]$, not on the specific distributions $\nu_i$ at hand). For the UCB algorithm,

$$\sup_{\substack{\nu_1, \dots, \nu_K \text{ with} \\ \text{supports in } [0,1]}} \bar{R}_T \leq O\left(\sqrt{TK\ln T}\right).$$

<u>Hint</u>: For small values of $\Delta_i$, the bound of the Proposition can be worse than the trivial $T$ bound...

<u>Proof</u> [ of the theorem based on the Proposition ]:

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

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

$$\mathbb{E}[y_t \mid I_t] = \mu_{I_t}$$

thus (by the tower rule) 
$$\mathbb{E}[y_t] = \mathbb{E}[\mu_{I_t}]$$
$$= \sum_{j} \mu_j \, \mathbb{E}\left[\mathbb{1}_{\{I_t = j\}}\right]$$

Summing over $t$: 
$$\mathbb{E}\left[\sum_{t=1}^{T} y_t\right] = \sum_{j=1}^{K} \mu_j \, \mathbb{E}[N_j(T)]$$

and ( in view of $T = \sum_{j} \mathbb{E}[N_j(T)]$ )

$$\bar{R}_T = \sum_{j} (\mu^* - \mu_j) \, \mathbb{E}[N_j(T)] = \sum_{j=1}^{K} \Delta_j \, \mathbb{E}[N_j(T)]$$
$$= \sum_{j:\,\Delta_j > 0} \Delta_j \, \mathbb{E}[N_j(T)]$$

it suffices to consider the suboptimal arms...

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

<u>NOTE</u>: Keep in mind the rewriting 
$$\bar{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^{T} y_t\right]$$
$$= \sum_{a=1}^{K} \Delta_a \, \mathbb{E}[N_a(T)]$$
as we will often use it!

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

→ It will show why this algorithm is called UCB:

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

estimate based on the raw performance
← exploitation 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 statements to the Hoeffding - Azuma inequality

replace $Y_1, Y_2, \dots$ by $1 - Y_1 \sim Y_2 \sim \dots$ then $\mu_a^{,} = 1 - \mu_a$ as well and $\Sigma_a^{,}$ supported on $[0,1]$

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

$\forall \delta \in (0,1)$,     $\mathbb{P}\left\{ \mu_j > \hat{\mu}_{j,t} - \sqrt{\dfrac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - t\delta$.

or

By symmetry:   $\forall \delta \in (0,1)$,     $\mathbb{P}\left\{ \mu_j < \hat{\mu}_{j,t} + \sqrt{\dfrac{\ln(1/\delta)}{2 N_j(t)}} \right\} \geq 1 - t\delta$

→ Application:   $N_i(T) = 1 + \sum\limits_{t=K+1}^{T} \mathbb{1}_{\{I_t = i\}}$

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

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

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

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

Indeed, we would otherwise have

$$\hat{\mu}_{a^*,\,t-1} + \sqrt{\frac{2\ln t}{N_{a^*}(t-1)}} \quad \geqslant \quad \mu^* \qquad\qquad \text{negation of (ii)}$$

$$= \mu_i + \Delta_i \qquad\qquad \text{definition of } \Delta_i$$

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

$$\geqslant \hat{\mu}_{i,\,t-1} + \sqrt{\frac{2\ln t}{N_i(t-1)}} \qquad \begin{array}{l}\text{negation}\\ \text{of (i)}\end{array}$$

the $\geqslant$ inequality between these two quantities would contradict $I_t = i$, that is, $i \in \arg\max_j \left\{ \hat{\mu}_{j,t} + \sqrt{2\ln t / N_j(t-1)} \right\}$

Thus, 
$$\mathbb{E}[N_i(T)] \leqslant 1 + \sum_{t=K+1}^{T} \mathbb{P}\left( \hat{\mu}_{i,\,t-1} > \mu_i + \sqrt{\frac{2\ln t}{N_i(t-1)}} \right) \qquad \begin{array}{l}\text{each} \leqslant t\delta\\ \text{where}\\ \delta = 2/t^4\end{array}$$

$$+ \sum_{t=K+1}^{T} \mathbb{P}\left( \hat{\mu}_{a^*,\,t-1} < \mu^* - \sqrt{\frac{2\ln t}{N_{a^*}(t-1)}} \right)$$

$$+ \mathbb{E}\left[ \sum_{t=K+1}^{T} \mathbb{1}_{\left\{ I_t = i \ \& \ N_i(t-1) \leqslant 8\ln t / \Delta_i^2 \right\}} \right] \qquad \begin{array}{l}8\ln t\\ \leqslant 8\ln T\end{array}$$

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

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{\text{deterministically upper bounded by}}$$

$$\leqslant 1 + 2\sum_{t=K+1}^{T} t^{-3} + \left( \frac{8\ln T}{\Delta_i^2} + 1 \right) - 1$$

as $\underbrace{I_t = i}_{}$ only if $N_i(t-1) \leqslant \dfrac{8\ln T}{\Delta_i^2}$

$$\leqslant 2 \int_1^{+\infty} t^{-3}\,dt$$

thus only if $N_i(t) \leqslant \dfrac{8\ln T}{\Delta_i^2} + 1$

$$= \left[ -t^{-2} \right]_1^{+\infty}$$

$$= 1$$

so that the total sum $\sum_{s=1}^{t} \mathbb{1}_{\{I_s = i\}} = N_i(t)$ is controlled $\forall t$ by this number

$-1$ because $I_1 = i$ is not included in the $\sum_{t=K+1}^{T}$

Thus:
$$\mathbb{E}[N_i(T)] \leqslant \frac{8\ln T}{\Delta_i^2} + 2$$

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

Let

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

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

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

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

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

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

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

$Y_t$ is drawn independently at random given $I_t$, thus by the very bandit model, this conditional expectation equals $\mu_{I_t}$

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

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

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

$\hookrightarrow$ in particular  $\mathbb{E}[M_t] \leqslant 1$  for all $t$

(2)  $\forall \varepsilon > 0, \ \forall \ell \geqslant 1,$    $\mathbb{P}\{Z_t \geqslant \varepsilon \ \text{and} \ N_a(t) = \ell\} \leqslant e^{-2\varepsilon^2/\ell}$

(3)  From these we will conclude.