

Convex functions & comparison to the best convex vector

↳ efficient forecaster (but worse bound...)

$X = \{(p_1, \dots, p_n) : \sum_j p_j = 1, \forall k p_k \geq 0\}$ is the simplex of convex weight vectors.

Setting (reminder):

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

1. The statistician and the opponent pick simultaneously $p_t \in X$ and a convex loss function $\ell_t: X \rightarrow [m, M]$

2. ℓ_t and p_t are revealed

↳ Regret: $R_T = \sum_{t=1}^T \ell_t(p_t) - \inf_{p \in X} \sum_{t=1}^T \ell_t(p)$ to be controlled in a uniform way

Fact: Convex functions are subdifferentiable on the interior of their domain of definition:

Let $f: \mathcal{D} \rightarrow \mathbb{R}$ be convex, where $\mathcal{D} \subseteq \mathbb{R}^N$ is convex:

$$\forall z \in \mathcal{D}, \exists \partial_{f,z} \in \mathbb{R}^N \mid \forall y \in \mathcal{D},$$

$$f(z) - f(y) \leq \partial_{f,z} \cdot (z - y)$$

$\partial_f(z) = \{ \text{the set of all possible such } \partial_{f,z} \}$ is called the subgradient of f at z

If f is differentiable at z , then $\partial_f(z) = \{ \nabla f(z) \}$.

Ex: $f(x) = |x|$
 $\partial_f(x) = \begin{cases} [-1, 1] & \text{if } x < 0 \\ \{1\} & \text{if } x > 0 \\ \{0\} & \text{if } x = 0 \end{cases}$

Application:

If $p_t \in X$ (ie $p_{jt} > 0 \forall j$) then $\exists \partial_{\ell_t}(p_t) \in \mathbb{R}^N \mid \forall p \in X, \ell_t(p_t) - \ell_t(p) \leq \partial_{\ell_t}(p_t) \cdot (p_t - p)$

Example:

Meta-statistical framework: $\ell_t(p) = \left(\sum_j p_j f_{jt} - y_t \right)^2$
 ℓ_t is differentiable, $\nabla \ell_t(p) = 2 \left(\sum_j p_j f_{jt} - y_t \right) \cdot (f_{jt})_j$
 ↳ Gradients in $[-G, G]^N$ with $G = 2B^2$ if $f_{jt}, y_t \in [a, b]$

Strategy:

Exponentiated Gradients (EG) with learning rate $\eta > 0$

$$p_t = \left(\frac{1}{N}, \dots, \frac{1}{N} \right) \text{ and}$$

$$p_{jt} = \exp\left(-\eta \sum_{s=1}^{t-1} (\partial_s(p_s))_j\right) / \sum_{k=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} (\partial_s(p_s))_k\right)$$
 → strategy easy to implement as soon as subgradients can be easily computed.

Theorem: Assume that the ℓ_t are picked such that

$$\forall t, \forall j, \forall p \in \mathcal{X}, \quad (\partial \ell_t(p))_j \in [-G, G]$$

Then the regret of EG tuned with $\eta > 0$ is controlled as

$$R_T = \sum_{t=1}^T \ell_t(p_t) - \inf_{p \in \mathcal{X}} \sum_{t=1}^T \ell_t(p) \leq \frac{\ln N}{\eta} + \eta \frac{G^2 T}{2}$$

In particular, the choice $\eta = \frac{1}{G} \sqrt{\frac{2 \ln N}{T}}$ leads to

$$R_T \leq G \sqrt{2 T \ln N}$$

Proof: Fix a p . Since by construction (cf. exponential weights), $p_t \in \mathcal{X}$ we have

$$\ell_t(p_t) - \ell_t(p) \leq \partial \ell_t(p_t) \cdot (p_t - p)$$

Thus

$$\begin{aligned}
 R_T &\leq \sup_{p \in \mathcal{X}} \sum_{t=1}^T \partial \ell_t(p_t) \cdot (p_t - p) \\
 &= \sup_{p \in \mathcal{X}} \sum_{t=1}^T \sum_{j=1}^N p_{jt} \tilde{\ell}_{jt} - \sum_{t=1}^T \sum_{j=1}^N p_j \tilde{\ell}_{jt}
 \end{aligned}$$

where we denoted by $\tilde{\ell}_{jt} = (\partial \ell_t(p_t))_j$ the components of the subgradients: we interpret them as pseudo-losses

By linearity in p of the last upper bound,

$$R_T \leq \sup_{p \in \mathcal{X}} \sum_{t=1}^T \sum_{j=1}^N p_{jt} \tilde{\ell}_{jt} - \sum_{t=1}^T \sum_j p_j \tilde{\ell}_{jt} = \sum_{t,j} p_{jt} \tilde{\ell}_{jt} - \min_K \sum_{t=1}^T \tilde{\ell}_{jt}$$

↳ We have performed a reduction to the linear case with vectors of losses

$(\tilde{\ell}_{1t} \dots \tilde{\ell}_{Nt})$ where $\tilde{\ell}_{jt} \in [-G, G]$

hence the stated bound.

Application: Oracle inequalities \rightarrow From individual sequences to stochastic sequences

Question: Let (Y_1, \dots, Y_T) be a sequence of iid random variables taking values in an arbitrary set \mathcal{Y} ; and $Q: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a bounded loss function, convex in its first argument, with bounded subgradients:

$$G = \sup_{y \in \mathcal{Y}} \sup_{p \in \mathcal{X}} \|\partial Q(p, y)\|_{\infty} < +\infty$$

Aim: Construct $\hat{\theta}_T = \hat{\theta}_T(Y_1, \dots, Y_T)$ such that

$$\mathbb{E}[Q(\hat{\theta}_T, Y)] \leq \inf_{\theta \in \mathcal{X}} \mathbb{E}[Q(\theta, Y)] + \varepsilon_T$$

where Y is independent of the Y_t with the same distribution and $\varepsilon_T \rightarrow 0$.

(The expectation \mathbb{E} is w.r.t. Y_1, \dots, Y_T and Y .)

Typical machine learning method:

Empirical risk minimization, possibly in a regularized way:

$$\hat{\theta}_T \in \operatorname{argmin}_{\theta \in \mathcal{X}} \left\{ \frac{1}{T} \sum_{t=1}^T Q(\theta, Y_t) + \lambda \operatorname{reg}(\theta) \right\}$$

with $\operatorname{reg}(\theta) = \|\theta\|_2$ or $\|\theta\|_1$ or ...
and λ to be tuned (called the regularization factor)

Our method: (as in our proof of Hoeffding's lemma)

1) Pretend data is sequential, while of course it is batch

$$\tilde{\theta}_1 = (1/N, \dots, 1/N) \quad \text{and by induction, with } \eta = \frac{1}{G} \sqrt{\frac{2 \ln N}{T}} :$$

$$\tilde{\theta}_{jt} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} (\partial Q(\tilde{\theta}_s, Y_s))_j\right)}{\sum_{k=1}^N \exp\left(-\eta \sum_{s=1}^{t-1} (\partial Q(\tilde{\theta}_s, Y_s))_k\right)}$$

2) Consider an average:

$$\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \tilde{\theta}_t$$

(note: we do not use y_t !)

Guarantee: $\epsilon_T = G \sqrt{\frac{2 \ln N}{T}}$

Proof: By the theorem on EG (exponentiated gradient) and given our choice of η :

$$\sum_{t=1}^T Q(\tilde{\theta}_t, y_t) - \inf_{\theta \in \mathcal{X}} \sum_{t=1}^T Q(\theta, y_t) \leq G \sqrt{2T \ln N}$$

Thus: $\forall \theta \in \mathcal{X}, \quad \frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t) \leq \frac{1}{T} \sum_{t=1}^T Q(\theta, y_t) + \epsilon_T$
 with $\epsilon_T = G \sqrt{\frac{2 \ln N}{T}}$

Now, since y_1, \dots, y_T and y are iid:

$$E\left[\frac{1}{T} \sum_{t=1}^T Q(\theta, y_t)\right] = E[Q(\theta, y)]$$

We conclude the proof by showing that $E\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t)\right] \geq E[Q(\hat{\theta}_T, y)]$

Indeed: $E\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t)\right] = E\left[\frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y)\right]$

$\tilde{\theta}_t$ only depends on y_1, \dots, y_{t-1} and $(y_1, \dots, y_{t-1}, y_t)$ and (y_1, \dots, y_{t-1}, y) have the same distribution
 $\geq Q\left(\frac{1}{T} \sum_{t=1}^T \tilde{\theta}_t, y\right)$
 by convexity and because the second argument is now the same for everyone!

We discard all intractability issues by using $E[\inf_{\theta \in \mathcal{X}} Q(\theta, y_t)]$ instead of considering each θ separately.

Extension: Oracle inequality for stationary data.

A stationary sequence (y_1, y_2, \dots) is by definition a sequence of random variables such that

$$\forall k \geq 1, \quad \forall t \geq 1, \quad (y_1, \dots, y_k) \text{ and } (y_{t+1}, \dots, y_{t+k}) \text{ have the same distribution.}$$

* Proof of performance? (\rightarrow explains where the construction comes from!)

What does EG guarantee?

$$\forall \theta \in \mathcal{X}, \quad \frac{1}{T} \sum_{t=1}^T Q(\Psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t) \leq \frac{1}{T} \sum_{t=1}^T Q(\theta, y_t) + \varepsilon_T$$

still $O(\sqrt{\ln T}/T)$

We take expectations in both sides to conclude :

- Stationarity implies in particular that all the y_t have the same distribution (but they are NOT independent in general); hence

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T Q(\theta, y_t) \right] = \mathbb{E} [Q(\theta, y_{T+1})]$$

- it also implies that $(\Psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t)$ and $(\Psi_t^{\text{EG}}(y_{T+2-t}, \dots, y_T), y_{T+1}) = (\hat{\theta}_t, y_{T+1})$ have the same distribution, so that

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T Q(\Psi_t^{\text{EG}}(y_1, \dots, y_{t-1}), y_t) \right] \\ &= \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T Q(\hat{\theta}_t, y_{T+1}) \right] \geq \mathbb{E} [Q(\hat{\theta}_T, y_{T+1})] \end{aligned}$$

↑
convexity of $\theta \mapsto Q(\theta, y_{T+1})$

- thus we proved, for our choice of $\hat{\theta}_T$:

$$\mathbb{E} [Q(\hat{\theta}_T, y_{T+1})] \leq \inf_{\theta \in \mathcal{X}} \mathbb{E} [Q(\theta, y_{T+1})] + \varepsilon_T$$

↑
same ε_T as for the iid case.

For a sequence of observations $z_1, \dots, z_t \in \mathcal{Y}$, EG outputs the following weights, which we re-define as functions of the information available:

$$\psi_1^{\text{EG}} = (\frac{1}{N}, \dots, \frac{1}{N})$$

and for $t \geq 2$:

$$\psi_t^{\text{EG}}(z_1, \dots, z_{t-1}) = \frac{\exp(-\eta \sum_{s=1}^{t-1} (\partial \ell(\psi_s^{\text{EG}}(z_1, \dots, z_{s-1}), z_s)))_j}{\sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} (\partial \ell(\psi_s^{\text{EG}}(z_1, \dots, z_{s-1}), z_s)))_k}$$

j -th component of
 $\psi_t^{\text{EG}}(z_1, \dots, z_{t-1}) \in \mathbb{R}^N$

! ψ_t^{EG} is an extremely complicated function of its arguments z_1, \dots, z_{t-1} (of the various calls to ψ_s^{EG} , with $s \leq t-1$, in its definition).

We define $\tilde{\theta}_t^{\text{EG}} = \psi_t^{\text{EG}}(y_{T+2-t}, \dots, y_T)$ for $t \geq 2$
 and $\tilde{\theta}_1^{\text{EG}} = (\frac{1}{N}, \dots, \frac{1}{N})$

We recommend: $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \tilde{\theta}_t^{\text{EG}}$

! Computing $\hat{\theta}_T$ is costly: the computation of each $\tilde{\theta}_t^{\text{EG}}$ requires $t-1$ steps, namely the computation of the

$$\psi_s^{\text{EG}}(y_{T+2-t}, \dots, y_{T-t+s}) \text{ for } 2 \leq s \leq t$$

no such intermediate quantity appears twice

\Rightarrow about $\sum_{t=1}^T (t-1) = O(\frac{T^2}{2})$ computation steps required.

Optimality of the $\sqrt{\frac{T}{2} \ln N}$ bound: asymptotic lower bound.

In the case of linear losses

By homogeneity we may assume that $l_{jt} \in [0,1]$ (ie, $m=0$ and $M=1$) for all t and j .

Theorem: All forecasting strategies of the statistician are such that

$$1 \leq \liminf_{N \rightarrow +\infty} \liminf_{T \rightarrow +\infty} \sup_{l_{jt} \in [0,1]} \frac{\sum_{t=1}^T p_{jt} l_{jt} - \min_k \sum_{t=1}^T l_{kt}}{\sqrt{\frac{T}{2} \ln N}}$$

the opponent does not need to react, it suffices to consider fixed-in-advance sequences (individual sequences)

Proof: We lower bound the sup by an expectation E_j assuming that the losses are given by random variables L_{jt} iid $\sim \text{Ber}(1/2)$

$\mathcal{F}_{t-1} = \sigma(L_{js}, s \leq t-1 \text{ and } j \in \{1, \dots, N\})$ is the information available at the beginning of round t : thus, p_t is \mathcal{F}_{t-1} -measurable

and by the tower rule,

by independence: $p_{jt} \perp L_{jt}$

$$E\left[\sum_j p_{jt} L_{jt}\right] = E\left[\sum_j p_{jt} E[L_{jt} | \mathcal{F}_{t-1}]\right] = 1/2 = E[L_{jt}] = 1/2$$

$$E[p_{jt} L_{jt}] = E[p_{jt}] E[L_{jt}] = 1/2$$

$$\sum_j E[p_{jt} L_{jt}] = 1/2$$

not so surprising: you cannot predict iid random variables!

* But * we will show that

$E\left[\min_k \sum_{t=1}^T L_{kt}\right]$ is smaller than $T/2$, even if "individually",

"crowd" effect (due to the central limit theorem) we have $\forall k, E\left[\sum_{t=1}^T L_{kt}\right] = T/2$

Let us summarize what we have so far:

$$\begin{aligned} \text{regret} &= \sup_{f_j \in \mathcal{Q}_1} \left[\sum_{j,t} p_{jt} \ell_{jt} - \min_k \sum_{t=1}^T \ell_{kt} \right] \geq \mathbb{E} \left[\sum_{j,t} p_{jt} \ell_{jt} - \min_k \sum_{t=1}^T \ell_{kt} \right] \\ &= \frac{T}{2} - \mathbb{E} \left[\min_k \sum_{t=1}^T \ell_{kt} \right] \\ &= \frac{\sqrt{T}}{2} \mathbb{E} \left[\max_{k=1..N} \frac{\sum_{t=1}^T (\frac{1}{2} - \ell_{kt})}{\frac{1}{2}\sqrt{T}} \right] \end{aligned}$$

Let that $\mathbb{E}[Z_{kT}] = 0$ and $\mathbb{E}[Z_{kT}] = 1$ independence } we denote $Z_{kT} = \frac{\sum_{t=1}^T (\frac{1}{2} - \ell_{kt})}{\frac{1}{2}\sqrt{T}}$

By the central limit theorem, $Z_{kT} \xrightarrow[T \rightarrow +\infty]{d} \mathcal{U}(\mathcal{Q}_1)$

We will show successively that

would be sufficient to prove $\lim_{T \rightarrow +\infty} \mathbb{E}[\max_k Z_{kT}] \geq \mathbb{E}[\max_k G_k]$

$$(1) \quad \mathbb{E} \left[\max_{k=1..N} Z_{kT} \right] \xrightarrow[T \rightarrow +\infty]{} \mathbb{E} \left[\max_{k \leq N} G_k \right]$$

where G_1, \dots, G_N iid $\sim \mathcal{U}(\mathcal{Q}_1)$

would be sufficient to prove $\lim_{N \rightarrow +\infty} \frac{\mathbb{E}[\max_{k \leq N} G_k]}{\sqrt{2 \ln N}} \geq 1$

$$(2) \quad \mathbb{E} \left[\max_{k \leq N} G_k \right] \underset{N \rightarrow +\infty}{\sim} \sqrt{2 \ln N}$$

which will conclude the proof of the theorem.

Point (1) relies on asymptotic uniform integrability.

Reminder: If $(Y_T)_{T \geq 1}$ is a sequence of \mathbb{R}^N -valued random variables with $Y_T \xrightarrow[T \rightarrow +\infty]{d} Y$, then for all continuous and bounded functions $f: \mathbb{R}^N \rightarrow \mathbb{R}$,

$$\mathbb{E}[f(Y_T)] \rightarrow \mathbb{E}[f(Y)]$$

Here: $Y_T = (Z_{kT})_{k \in \{1, \dots, N\}} \in \mathbb{R}^N$

and $Y_T \xrightarrow{d} \begin{bmatrix} G_1 \\ \vdots \\ G_N \end{bmatrix} \sim \mathcal{U}(0, I_N)$

by independence of the components

Note: in general $Y_T \xrightarrow{d} Y$ does not imply $f(Y_T) \rightarrow f(Y)$ if f is just assumed continuous.
Ex: $Y_T = T^2 \mathbb{1}_{\{U_T \in [0, 1/T]\}}$ where U_T iid $\sim \mathcal{U}(\mathcal{Q}_1)$. We have $Y_T \rightarrow 0$ in probability but $\mathbb{E} Y_T = T \rightarrow +\infty$

* but * $f(z_1, \dots, z_N) = \max_{k \leq N} z_k$ is not bounded (though being continuous).

Definition: $(f(Y_T))_{T \geq 1}$ is a.u.i. (asymptotically uniformly integrable)

if
$$\lim_{L \rightarrow +\infty} \limsup_{T \rightarrow +\infty} E[|f(Y_T)| \mathbb{1}_{\{|f(Y_T)| > L\}}] = 0$$

Lemma: If $\left\{ \begin{array}{l} Y_T \xrightarrow{d} Y \\ (f(Y_T)) \text{ is a.u.i.} \\ f \text{ continuous} \end{array} \right\}$ then (a) $f(Y) \in L^1$ and $f(Y_T) \xrightarrow{d} f(Y)$ and (b) $E[f(Y_T)] \xrightarrow{T \rightarrow +\infty} E[f(Y)]$

(a) is added to page in this document

Proof: \hookrightarrow (b) We may assume with no loss of generality that $f \geq 0$
 (otherwise, we consider $f^+ = \max\{f, 0\}$ and $f^- = \max\{-f, 0\}$ the positive and negative parts of f , and use $f = f^+ - f^-$ and the fact that $f(Y) \in L^1$)

$$\begin{aligned} & \therefore |E[f(Y)] - E[f(Y_T)]| \\ & \leq |E[f(Y)] - E[\min\{f(Y), L\}]| \\ & \quad + |E[\min\{f(Y), L\}] - E[\min\{f(Y_T), L\}]| \\ & \quad + |E[\min\{f(Y_T), L\}] - E[f(Y_T)]| \end{aligned}$$

\rightsquigarrow tends to 0 as $T \rightarrow +\infty$, since $\min\{f(\cdot), L\}$ is continuous, bounded in $[0, L]$; and $Y_T \xrightarrow{d} Y$

thus, $\forall L$:
$$\begin{aligned} & \limsup_{T \rightarrow +\infty} |E[f(Y)] - E[f(Y_T)]| \\ & \leq |E[f(Y)] - E[\min\{f(Y), L\}]| \\ & \quad + \limsup_{T \rightarrow +\infty} |E[\min\{f(Y_T), L\}] - E[f(Y_T)]| \end{aligned}$$

taking $\limsup_{L \rightarrow +\infty}$ in both sides (the left-hand side is independent of L)

and using that $\limsup_{T \rightarrow +\infty} (a_T + b_T) \leq \limsup_{T \rightarrow +\infty} a_T + \limsup_{T \rightarrow +\infty} b_T$
 think in terms of largest adherence values]

we get
$$\begin{aligned} \limsup_{T \rightarrow +\infty} |E[f(Y)] - E[f(Y_T)]| & \leq \limsup_{L \rightarrow +\infty} |E[f(Y)] - E[\min\{f(Y), L\}]| \\ & \quad + \limsup_L \limsup_T |E[\min\{f(Y_T), L\}] - E[f(Y_T)]| \end{aligned}$$

By monotone convergence,
(or dominated)

$$E[\min\{f(Y), L\}] \xrightarrow{L \rightarrow +\infty} E[f(Y)]$$

And by the u.a.i property,

$$\limsup_{L \rightarrow +\infty}$$

$$\limsup_{T \rightarrow +\infty}$$

$$\begin{aligned} & \left| E[\min\{f(Y_T), L\}] - E[f(Y_T)] \right| = 0 \\ & = \left| E[(f(Y_T) - L) \mathbb{1}_{\{f(Y_T) > L\}}] \right| \\ & \stackrel{\text{because } f \text{ assumed } \geq 0}{\leq} E[f(Y_T) \mathbb{1}_{\{f(Y_T) > L\}}] \end{aligned}$$

Application:

Need only to check that $\max_{k \leq N} Z_{kT}$ is u.a.i. to get (i) from the lemma.

Indeed, given $L > 0$,

$$E\left[\left| \max_{k \leq N} Z_{kT} \right| \mathbb{1}_{\left\{ \left| \max_{k \leq N} Z_{kT} \right| > L \right\}} \right]$$

$$= \int_0^{+\infty} P\left(\left| \max_{k \leq N} Z_{kT} \right| \mathbb{1}_{\left\{ \left| \max_{k \leq N} Z_{kT} \right| > L \right\}} \geq u \right) du$$

a well-known fact: if $X \geq 0$,
 $E[X] = \int_0^{+\infty} P(X \geq u) du$

and, by Fubini-Tonelli:

$$\begin{aligned} E[X] &= \int_{[0, +\infty)} x dP^X(x) = \int_{[0, +\infty)}^2 \mathbb{1}_{\{0 \leq u \leq x\}} dP^X(x) du \\ &= \int_{[0, +\infty)} P(X \geq u) du \end{aligned}$$

equals

$$\begin{cases} P\left(\left| \max_{k \leq N} Z_{kT} \right| > L \right) & \text{if } u \leq L \\ P\left(\left| \max_{k \leq N} Z_{kT} \right| \geq u \right) & \text{if } u > L \end{cases}$$

But, for all $u > 0$:

$$\begin{aligned} P\left(\left| \max_{k \leq N} Z_{kT} \right| \geq u \right) &\leq P\left(\exists j \leq N \mid |Z_{jT}| \geq u \right) && \text{union bound + same distribution} \\ &\leq N \times P\left(|Z_{1T}| \geq u \right) \\ &\leq N \times \frac{\text{Var } Z_{1T}}{u^2} = \frac{N}{u^2} && \begin{matrix} E Z_{1T} = 0 \\ \text{and Chebyshev} \\ \text{if: Var } Z_{1T} = 1 \end{matrix} \end{aligned}$$

Substituting above, the conclusion is obtained:

$$\begin{aligned} E\left[\left| \max_{k \leq N} Z_{kT} \right| \mathbb{1}_{\left\{ \left| \max_{k \leq N} Z_{kT} \right| > L \right\}} \right] &\leq \int_0^L \frac{N}{L^2} du + \int_L^{+\infty} \frac{N}{u^2} du = \frac{2N}{L} \\ &\text{bound independent of } T \text{ and such that } \rightarrow 0 \text{ as } L \rightarrow +\infty \end{aligned}$$

Alternative way:

$$\begin{aligned}
 & \mathbb{E} \left[\left| \max_{k \leq N} Z_{kT} \right| \mathbb{1}_{\left\{ \left| \max_{k \leq N} Z_{kT} \right| > L \right\}} \right] \\
 \stackrel{\text{Cauchy-Schwarz}}{\leq} & \underbrace{\left(\mathbb{E} \left[\max_{k \leq N} Z_{kT}^2 \right] \right)^{1/2}}_{\leq \left(\sum_{k \leq N} \mathbb{E} Z_{kT}^2 \right)^{1/2} = \left(N \operatorname{Var}(Z_{1T}) \right)^{1/2} = N^{1/2}} \underbrace{\left(\mathbb{P} \left\{ \left| \max_{k \leq N} Z_{kT} \right| > L \right\} \right)^{1/2}}_{\leq \sqrt{N}/L \text{ as proved above}} \leq N/L \rightarrow 0
 \end{aligned}$$

Note (credits: Clement Borel + Cherlin Ga.)

follows from $\mathcal{Y}_T \subset \mathcal{Y}$
 f continuous

Proof of property (a): $(f(\mathcal{Y}_T))$ u.a.i and $f(\mathcal{Y}_T) \rightarrow f(\mathcal{Y})$
 implies $f(\mathcal{Y}) \in \mathbb{L}^1$ (and $f(\mathcal{Y}_T) \in \mathbb{L}^1$ $\forall T$ sufficiently large)

Indeed: By u.a.i, $\exists L_0$ | $\limsup_{T \rightarrow +\infty} E[|f(\mathcal{Y}_T)| \mathbb{1}_{|f(\mathcal{Y}_T)| > L_0}] < +\infty$

in particular, $\limsup_{T \rightarrow +\infty} E[|f(\mathcal{Y}_T)|] < +\infty$

Thus, $(E[|f(\mathcal{Y}_T)|])_{T \geq 1}$ is bounded

By Skorokhod, we can assume with no loss of generality that

$f(\mathcal{Y}_T) \rightarrow f(\mathcal{Y})$ a.s.

so that Fatou's lemma indicates $E[|f(\mathcal{Y})|] \leq \liminf_{T \rightarrow +\infty} E[|f(\mathcal{Y}_T)|] < +\infty$

It only remains to prove (2):

$$\mathbb{E} \left[\max_{k \leq N} G_k \right] \sim \sqrt{2 \ln N}$$

where $G_j \text{ iid} \sim \mathcal{U}(0,1)$

* Easy (and useless part for our proof): $\mathbb{E} \left[\max_{k \leq N} G_k \right] \leq \sqrt{2 \ln N}$

Lemma: If V_1, \dots, V_N are such that $\mathbb{E}[e^{dV_j}] \leq e^{d^2\sigma^2/2} \forall j, d > 0$
then $\mathbb{E} \left[\max_{k \leq N} V_k \right] \leq \sqrt{2\sigma^2 \ln N}$

Note: Valid even without independence \rightarrow independence will be needed for the lower bound

If $G_k \sim \mathcal{U}(0, \sigma^2)$ then $\mathbb{E}[e^{dG_k}] = e^{d^2\sigma^2/2}$

If $G_k \in [m, M]$ with $\mathbb{E}[G_k] = 0$ then (Hoeffding's lemma) $\mathbb{E}[e^{dG_k}] \leq e^{d^2(M-m)^2/8}$

\rightarrow We call subgaussian the random variables V_k satisfying the assumptions of the lemma.

Proof:

Pisier's argument:

$$\begin{aligned} \mathbb{E} \left[\max_{k \leq N} V_k \right] &\stackrel{\text{Jensen's ineq. } d > 0}{\leq} \frac{1}{d} \ln \mathbb{E} \left[e^{d \max_{k \leq N} V_k} \right] \\ &\leq \frac{1}{d} \ln \mathbb{E} \left[\sum_{k=1}^N e^{dV_k} \right] \leq \frac{1}{d} \ln \left(N e^{d^2\sigma^2/2} \right) \\ &= \frac{\ln N}{d} + d \frac{\sigma^2}{2} = \sqrt{2\sigma^2 \ln N} \end{aligned}$$

choose $d^* = \sqrt{\frac{2 \ln N}{\sigma^2}}$

* « Difficult » part of the proof:

$$\liminf_{N \rightarrow +\infty} \frac{\mathbb{E} \left[\max_{k \leq N} G_k \right]}{\sqrt{2 \ln N}} \geq 1$$

(Credit: Pascal Nassart explained to me how to proceed efficiently!)

We use: let $F_N = \Phi^N$ be the cumulative distribution function of

$\max_{k \leq N} G_k$, where Φ is the cumulative distribution function of the $\mathcal{U}(0,1)$ distribution.

F_N is a bijection $\mathbb{R} \rightarrow (0,1)$ so that, for $U \sim \mathcal{U}_{[0,1]}$,
 $F_N^{-1}(U) = \Phi^{-1}(U^{1/N})$ and $\max_{k \leq N} G_k$ have the same distribution.

in particular,

$$E[\max_{k \leq N} G_k] = E[\Phi^{-1}(U^{1/N})] = \int_0^1 \Phi^{-1}(u^{1/N}) du$$

$$= \int_0^{1/2^N} \underbrace{\Phi^{-1}(u^{1/N})}_{\leq 0} du + \int_{1/2^N}^1 \underbrace{\Phi^{-1}(u^{1/N})}_{\geq 0} du$$

we separate the integrals depending on the sign of the integrand.

is ≤ 0 and is ≥ 0

$$E[\max_{k \leq N} G_k \mathbb{1}_{\{\max_{k \leq N} G_k \leq 0\}}]$$

$$= E[G_1 \mathbb{1}_{\{G_1 \leq 0\}} \mathbb{1}_{\{\max_{k \leq N} G_k \leq 0\}}]$$

$$\geq E[G_1 \mathbb{1}_{\{G_1 \leq 0\}}]$$

fix any $\delta > 0$, for N large enough, $\delta \geq 1/2^N$

$$\geq \int_{\delta}^1 \Phi^{-1}(u^{1/N}) du \geq (1-\delta) \Phi^{-1}(\delta^{1/N})$$

Φ^{-1} increasing

is

$$|E[\max_{k \leq N} G_k \mathbb{1}_{\{\max_{k \leq N} G_k > 0\}}]|$$

$$\leq |E[G_1 \mathbb{1}_{\{G_1 > 0\}}]|$$

$$\leq \sqrt{E G_1^2} \sqrt{P\{G_1 > 0\}} = 1/\sqrt{2}$$

Cauchy Schwarz

Thus we have proved so far:

$$\liminf_{N \rightarrow \infty} \frac{E[\max_{k \leq N} G_k]}{\sqrt{2 \ln N}} \geq \sup_{\delta \in (0,1)} \liminf_{N \rightarrow \infty} \frac{(1-\delta) \Phi^{-1}(\delta^{1/N})}{\sqrt{2 \ln N}}$$

It's only a matter of real analysis / calculus now:

We will successfully prove

$$(a) \quad 1 - \Phi(x) \underset{x \rightarrow +\infty}{\sim} \frac{e^{-x^2/2}}{x\sqrt{2\pi}}$$

$$(b) \quad \ln \frac{1}{1-u} \sim \left(\Phi^{-1}(u)\right)^2/2 \quad \text{as } u \nearrow 1$$

(c) thus the desired result.

roof of (a): by two integration by parts (IBP)

$$\begin{aligned} 1 - \Phi(x) &= \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} e^{-t^2/2} dt = \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} t e^{-t^2/2} \cdot \frac{1}{t} dt \\ &\stackrel{\text{first IBP}}{=} \underbrace{\left[-\frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t} \right]_x^{+\infty}}_{= \frac{e^{-x^2/2}}{x\sqrt{2\pi}}} - \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \frac{e^{-t^2/2}}{t^2} dt \end{aligned}$$

Second IBP: $0 \leq \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} t e^{-t^2/2} \frac{1}{t^3} dt$

$$= \left[-\frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t^3} \right]_x^{+\infty} - \int_x^{+\infty} \frac{e^{-t^2/2}}{t^4} dt$$

thus $\leq \frac{1}{\sqrt{2\pi}} \frac{e^{-x^2/2}}{x^3}$ $\underbrace{\quad}_{>0}$

so that the $\frac{e^{-x^2/2}}{x\sqrt{2\pi}}$ term was dominating in what the first IBP yielded, hence the result.

roof of (b):

Recall that if $f(x) \sim g(x)$ and $\ln f(x)$ is bounded away from 0 in the limit, then $\ln f(x) \sim \ln g(x)$

$$\text{indeed: } \frac{\ln g(x)}{\ln f(x)} = 1 + \frac{\ln(g(x)/f(x))}{\ln f(x)} \rightarrow 1$$

In particular

$$\ln \frac{1}{1-\Phi(x)} \underset{x \rightarrow +\infty}{\sim} \ln \left(e^{x^2/2} x \sqrt{2\pi} \right) \underset{x \rightarrow +\infty}{\sim} \frac{x^2}{2}$$

with $x = \Phi^{-1}(v) \rightarrow +\infty$ as $v \rightarrow 1$ (by composing limits):

$$\ln \frac{1}{1-v} \underset{v \rightarrow 1}{\sim} \frac{(\Phi^{-1}(v))^2}{2}$$

$$(c) \quad \Phi^{-1}(\delta^{1/N}) \sim \sqrt{2 \ln \frac{1}{1-\delta^{1/N}}} \\ \sim \sqrt{2 \ln N} \quad \text{as } N \rightarrow +\infty$$

$$\text{given } \delta^{1/N} \xrightarrow[N \rightarrow +\infty]{} 1$$

$$\text{given } 1 - \delta^{1/N} = 1 - \exp\left(\frac{1}{N} \ln \delta\right) \\ \sim -\frac{1}{N} \ln \delta = \frac{1}{N} \ln \frac{1}{\delta}$$

thus:

$$\liminf_N \frac{\mathbb{E}[\max_{k \leq N} G_k]}{\sqrt{2 \ln N}} \geq \sup_{\delta \in (0,1)} \left\{ (1-\delta) \times \underbrace{\liminf_{N \rightarrow +\infty} \frac{\Phi^{-1}(\delta^{1/N})}{\sqrt{2 \ln N}}}_{=1} \right\} = 1$$

ExerciseFinite-time lower bound(also known as:
Non-asymptotic lower bound)

It's a difficult exercise, try to do as much as possible!
Part of the difficulty comes from the manipulation of Kullback-Leibler divergence, we will study their properties in detail in a few weeks when studying the lower bounds for stochastic bandits.

Let \mathbb{P}_0 be a probability such that the losses L_{kt} are i.i.d. $\sim \text{Ber}(\frac{1}{2})$ $\forall k, t$

For $j \in \{1, \dots, N\}$ let \mathbb{P}_j be a probability such that:

- all losses L_{kt} are independent
- $L_{kt} \sim \text{Ber}(\frac{1}{2})$ $\forall t, \forall k \neq j$
- $L_{jt} \sim \text{Ber}(\frac{1}{2} - \epsilon)$ $\forall t$

(1) Show that

$$\sup_{\{p_{jt} \in [0,1]\}} \left\{ \sum_{j,t} p_{jt} L_{jt} - \min_k \sum_{t=1}^T L_{kt} \right\} \geq T \epsilon \max_{k=1, \dots, N} \left\{ 1 - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_k [p_{kt}] \right\}$$

REMINDER or QUICK INTRODUCTION TO A NEW OBJECT:

Recall that the Kullback-Leibler divergence between two probability distributions μ and ν equals:

$$KL(\nu, \mu) = \begin{cases} +\infty & \text{if } \mu \text{ not absolutely continuous w.r.t } \nu \\ \int \left(\frac{d\nu}{d\mu} \ln \frac{d\nu}{d\mu} \right) d\mu & \in [0, +\infty] \\ = \int \left(\ln \frac{d\nu}{d\mu} \right) d\nu & \text{if } \nu \ll \mu \end{cases}$$

(KL well-defined as $x \ln x$ is $\geq -1/e$ over $(0, +\infty)$)

KL has the following properties:

- $KL(\nu \otimes \nu', \mu \otimes \mu') = KL(\nu, \mu) + KL(\nu', \mu')$ \leftarrow by Fubini-Tonelli
- \forall random variables X, Y , $KL(\nu^X, \mu^X) \leq KL(\nu, \mu)$ \leftarrow data compression inequality (we take it for granted and will prove it later)

We denote by $KL(p, q) = KL(\text{Ber}(p), \text{Ber}(q))$

$$= p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q}$$

the KL-divergence between two Bernoulli distributions.

2) Prove Fano's lemma as revisited by Lucien Birgé:

let (Ω, \mathcal{F}) be a measurable space

let Q_1, \dots, Q_N be probability distributions over Ω

let (A_1, \dots, A_N) be a partition of Ω

then

$$\min_{j=1, \dots, N} Q_j(A_j) \leq \max \left\{ \frac{2e}{2e+1}, \frac{\bar{K}}{\ln N} \right\}$$

where $\bar{K} = \frac{1}{N-1} \sum_{j=2}^N KL(Q_j, Q_1)$

Consider $p = \frac{1}{N-1} \sum_{j=2}^N Q_j(A_j)$ and $q = \frac{1}{N-1} \sum_{j=2}^N Q_1(A_j)$

and denote $a = \min_{j=1, \dots, N} Q_j(A_j)$. Show that:

- $kl(p, q) \leq \bar{K}$ / to that end, first prove that KL is jointly convex (it is a consequence of the data processing inequality, explain why)

- $kl(p, q) / \ln N \geq a$ for $a \geq \frac{2e}{2e+1}$

(3) Extend Fano's lemma to random variables so as to show

$$\min_{k=1, \dots, N} \mathbb{E}_k \left[\frac{1}{T} \sum_{t=1}^T P_{kt} \right] \leq \max \left\{ \frac{2e}{2e+1}, \frac{\bar{K}'}{\ln N} \right\}$$

where $\bar{K}' = \frac{1}{N-1} \sum_{j=2}^N KL(\mathbb{P}_j^L, \mathbb{P}_1^L)$ with $L = (L_{kt})_{\substack{t \leq T \\ k \leq N}}$

(4) Show that $\bar{K}' \leq 5T\varepsilon^2$ for $\varepsilon \leq 1/10$.

(5) Conclude to the following theorem:

rem: For all strategies, for all $N \geq 2$ and $T \geq 17 \ln N$,
 $\sup R_T \geq 0.06 \sqrt{T \ln N}$