

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  $\sup_{l_{jt}}$  by an expectation  $\mathbb{E}_j$  assuming that the losses are given by random variables  $L_{jt} \text{ iid} \sim \text{Ber}(1/2)$

$\mathcal{F}_{t-1} = \sigma(L_{js}, s \leq t-1 \text{ and } j \in \{1, \dots, N\})$  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:  $p_{jt} \perp\!\!\!\perp L_{jt}$

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

$$= \mathbb{E}[L_{jt}] = \frac{1}{2}$$

$$\begin{aligned} p_{jt} L_{jt} &= \mathbb{E}[p_{jt}] \mathbb{E}[L_{jt}] \\ &= \mathbb{E}[p_{jt}] \frac{1}{2} \\ \sum_j \mathbb{E}[p_{jt} L_{jt}] &= \frac{1}{2} \end{aligned}$$

not so surprising:  
you cannot predict iid random variables!

\* But \*  
we will show that

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

we have

$$\mathbb{E} \left[ \sum_{t=1}^T L_{kt} \right] = T/2$$

"crowd"  
effect

(due to the central limit theorem)

Let us summarize what we have so far:

rest-case  

$$\text{regret} = \sup_{f_j \in [0,1]} \left[ \sum_{j,t} p_j f_j - \min_k \sum_{t=1}^T L_{kt} \right] \geq \mathbb{E} \left[ \sum_{j,t} p_j f_j - \min_k \sum_{t=1}^T L_{kt} \right]$$

$$= \frac{T}{2} - \mathbb{E} \left[ \min_k \sum_{t=1}^T L_{kt} \right]$$

$$= \frac{\sqrt{T}}{2} \mathbb{E} \left[ \max_{k=1..N} \frac{\sum_{t=1}^T (\frac{1}{2} - L_{kt})}{\frac{1}{2}\sqrt{T}} \right]$$

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

By the central limit theorem,  $Z_{kT} \xrightarrow[T \rightarrow +\infty]{d} \mathcal{U}(0,1)$

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)  $\mathbb{E}[\max_{k=1..N} Z_{kT}] \xrightarrow[T \rightarrow +\infty]{} \mathbb{E}[\max_{k \leq N} G_k]$   
 where  $G_1, \dots, G_N$  iid  $\sim \mathcal{U}(0,1)$

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)  $\mathbb{E}[\max_{k \leq N} G_k] \sim \sqrt{2 \ln N}$   
 as  $N \rightarrow +\infty$

which will conclude the proof of the theorem.

Point (1) relies on asymptotic uniform integrability.

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

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:

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

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

by independence of the components

\* 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)

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

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

of (a) is added to page 10 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  $f(y) \in L^1$ )

$$\begin{aligned} \forall L: & \quad |\mathbb{E}[f(y)] - \mathbb{E}[f(y_T)]| \\ & \leq |\mathbb{E}[f(y)] - \mathbb{E}[\min\{f(y), L\}]| \\ & \quad + |\mathbb{E}[\min\{f(y), L\}] - \mathbb{E}[\min\{f(y_T), L\}]| \\ & \quad + |\mathbb{E}[\min\{f(y_T), L\}] - \mathbb{E}[f(y_T)]| \end{aligned}$$

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

$$\begin{aligned} \text{Thus, } \forall L: & \quad \limsup_{T \rightarrow +\infty} |\mathbb{E}[f(y)] - \mathbb{E}[f(y_T)]| \\ & \leq |\mathbb{E}[f(y)] - \mathbb{E}[\min\{f(y), L\}]| \\ & \quad + \limsup_{T \rightarrow +\infty} |\mathbb{E}[\min\{f(y_T), L\}] - \mathbb{E}[f(y_T)]| \end{aligned}$$

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

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

By monotone convergence,  
(or dominated)

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

And by the u.a.i property,  
 $\limsup_{L \rightarrow +\infty}$   $\limsup_{T \rightarrow +\infty}$

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

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

Indeed, given  $L > 0$ ,

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

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

and, by Fubini-Tonelli:

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

equals

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

But, for all  $u > 0$ :

$$\begin{aligned} P(|\max_{k \leq N} Z_{kT}| \geq u) &\leq P(\exists j \leq N \mid |Z_{jT}| \geq u) \\ &\leq N \times P(|Z_{1T}| \geq u) \\ &\leq N \times \frac{\text{Var } Z_{1T}}{u^2} = \frac{N}{u^2} \end{aligned}$$

union bound + same distribution  
 $E Z_{1T} = 0$  and Chebyshev  
if  $\text{Var } Z_{1T} = 1$

Substituting above, the conclusion is obtained:

$$E[|\max_{k \leq N} Z_{kT}| \mathbb{1}_{\{|\max_{k \leq N} Z_{kT}| > L\}}] \leq \int_0^L \frac{N}{L^2} du + \int_L^{+\infty} \frac{N}{u^2} du = \frac{2N}{L}$$

bound independent of  $T$   
and such that  $\rightarrow 0$  as  $L \rightarrow +\infty$

Note (credits: Clement Berenfeld  
+ Chenlin Gu.)

follows from  $Y_T \subset Y$   
&  $f$  continuous

Proof of property (a):  $(f(Y_T))$  u.a.i. and  $f(Y_T) \rightarrow f(Y)$   
implies  $f(Y) \in L^1$  (and  $f(Y_T) \in L^1 \forall T$ )

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

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

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

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

$f(Y_T) \rightarrow f(Y)$  a.s.

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



It only remains to prove (2):

$$\mathbb{E} \left[ \max_{k \leq N} G_k \right] \sim \sqrt{2 \ln N} \quad \text{where } G_j \text{ iid} \sim \mathcal{U}(0,1)$$

\* 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)  $\mathbb{E}[e^{dG_k}] \leq e^{d^2 (M-m)^2 / 8}$

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

Proof:

Pisier's argument:

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

\* « 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  $\Phi$  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,

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

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

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

is  $\leq 0$  and is  $\geq 0$

$$\begin{aligned} & \mathbb{E}[\max_{k \leq N} G_k \mathbb{1}_{\{\max_{k \leq N} G_k \leq 0\}}] \\ & \geq \mathbb{E}[G_1 \mathbb{1}_{\{\max_{k \leq N} G_k \leq 0\}}] \\ & = \mathbb{E}[G_1 \mathbb{1}_{\{G_1 \leq 0\}} \mathbb{1}_{\{\max_{k \leq N} G_k \leq 0\}}] \\ & \geq \mathbb{E}[G_1 \mathbb{1}_{\{G_1 \leq 0\}}] \end{aligned}$$

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

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

$\Phi^{-1}$  increasing

is

$$\begin{aligned} & |\mathbb{E}[\max_{k \leq N} G_k \mathbb{1}_{\{\max_{k \leq N} G_k > 0\}}]| \\ & \leq |\mathbb{E}[G_1 \mathbb{1}_{\{G_1 > 0\}}]| \\ & \leq \sqrt{\mathbb{E}G_1^2} \sqrt{\mathbb{P}\{G_1 > 0\}} = 1/\sqrt{2} \end{aligned}$$

Cauchy Schwarz

Thus we have proved so far:

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

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

We will successfully prove

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

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

(c) thus the desired result.

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

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

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

$$= \left[ -\frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t^3} \right]_x^{+\infty} - \int_x^{+\infty} \frac{e^{-t^2/2}}{t^4} dt \geq 0$$

thus  $\leq \frac{1}{\sqrt{2\pi}} \frac{e^{-x^2/2}}{x^3}$

so that the  $e^{-x^2/2}/x\sqrt{2\pi}$  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  $\ln f(x) \sim \ln g(x)$

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



In particular

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

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

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

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

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

$$\sim \sqrt{2 \ln N} \quad \text{as } N \rightarrow +\infty$$

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

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 divergences, 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
- $L_{kt} \sim \text{Ber}(\frac{1}{2})$   $\forall t, \forall k \neq j$
- $L_{jt} \sim \text{Ber}(\frac{1}{2} - \varepsilon)$   $\forall t$

$$(1) \text{ Show that } \sup_{\ell_{jt} \in \{0,1\}} \left\{ \sum_{jt} p_{jt} \ell_{jt} - \min_k \sum_{t=1}^T \ell_{kt} \right\} \\ \geq T \varepsilon \max_{k=1, \dots, N} \left\{ 1 - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_k [p_{kt}] \right\}$$

REMINDER or QUICK INTRODUCTION TO A NEW OBJECT:

Recall that the Kullback-Leibler divergence between two probability distributions  $\mu$  and  $\nu$  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 variable  $X, KL(\nu^X, \mu^X) \leq KL(\nu, \mu)$   $\leftarrow$  data compression inequality (we take it for granted and will prove it later)

$$\text{We denote by } \mathcal{KL}(p, q) = KL(\text{Ber}(p), \text{Ber}(q)) \\ = p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} \quad \text{the KL-divergence between two Bernoulli distributions.}$$

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

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

let  $Q_1, \dots, Q_N$  be probability distributions over  $\Omega$

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

then

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

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

Hints:

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  $a \geq \frac{2e}{2e+1}$

$KL(E_{\sum_{t=1}^T Q_t}, E_{\sum_{t=1}^T Q_t}) \leq KL(P, Q)$   
 for any random variable  $Z_t$  on  $(\Omega, \mathcal{F})$  equipped with  $P$  or  $Q$

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

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

where  $\bar{K}' = \frac{1}{N-1} \sum_{j=2}^N KL(P_j^L, P_1^L)$  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:

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

## 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.  $\ell_t$  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:

$$\forall x \in \mathring{\mathcal{D}}, \exists \partial f_x \in \mathbb{R}^N \mid \forall y \in \mathcal{D},$$

$$f(x) - f(y) \leq \partial f_x \cdot (x - y)$$

$\partial f(x) = \{\text{the set of all possible such } \partial f_x\}$  is called the subgradient of  $f$  at  $x$

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

Ex:  $f(x) = |x|$

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

Application:

If  $p_t \in X$  (i.e.  $p_{jt} \geq 0 \forall j$ ) 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:

$$\ell_t(p) = \left( \sum_j p_j f_{jt} - y_t \right)^2$$

$\ell_t$  is differentiable,

$$\nabla \ell_t(p) = 2 \left( \sum_j p_j f_{jt} - y_t \right) \cdot f_{jt}$$

↳ Gradients in  $[-G, G]^N$  with  $G = 2B^2$  if  $f_{jt}, y_t \in [0, B]$

Strategy:

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

$$p_1 = \left( \frac{1}{N}, \dots, \frac{1}{N} \right) \text{ and}$$

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

Theorem: Assume that the  $\ell_t$  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,

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

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

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

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

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

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

(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  $\lambda$  to be tuned (called the regularization factor)

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

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

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

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



2) Consider an average:

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

(note: we do not use  $\theta_t$ !)

Guarantee:

$$\mathcal{E}_T = G \sqrt{\frac{2 \ln N}{T}}$$

Proof:

By the theorem on EG (exponentiated gradient) and given our choice of  $\eta$ :

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

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

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

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

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

Indeed:  $\mathbb{E} \left[ \frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y_t) \right] = \mathbb{E} \left[ \frac{1}{T} \sum_{t=1}^T Q(\tilde{\theta}_t, y) \right]$   
 $\uparrow$   
 $\tilde{\theta}_t$  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  
 $\geq \mathbb{E} \left[ Q \left( \frac{1}{T} \sum_{t=1}^T \tilde{\theta}_t, y \right) \right]$   
 $\uparrow$   
 by convexity and  $Q$  because the second argument is now the same for everyone!

Extension: Oracle inequality for stationary data.

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

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

Construct  $\hat{\theta}_T = \hat{\theta}_T(y_1, \dots, y_T)$  such that, under the same conditions on  $Q$  as above:

$$E[Q(\hat{\theta}_T, y_{T+1})] \leq \inf_{\theta} E[Q(\theta, y_{T+1})] + G \sqrt{\frac{2 \ln N}{T}}.$$

Hint: the sequential ordering of the  $y_t$  may be different from the "natural" order used in the iid case. Try to mimic the key equality in distribution used in the iid case, that is, the fact that  $(\tilde{\theta}_t, y_t)$  and  $(\tilde{\theta}_t, y_{T+1})$  have the same distribution for some well-chosen  $\tilde{\theta}_t$ .

Solution: see next pages; not very clear (I hope that I ended up being clearer than that in class...).

## Exercise / Oracle inequality for stationary data.

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:

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

and for  $t \geq 2$ :

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

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

⚠  $\psi_t^{\text{EG}}$  is an extremely complicated function of its arguments  $z_1, \dots, z_{t-1}$  (cf the various calls to  $\psi_s^{\text{EG}}$ , 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  $\tilde{\Theta}_1^{\text{EG}} = (1/N, \dots, 1/N)$

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

⚠ Computing  $\hat{\Theta}_T^{\text{EG}}$  is costly: the computation of each  $\tilde{\Theta}_t^{\text{EG}}$  requires  $t-1$  steps, namely the computation of the  $\psi_s^{\text{EG}}(y_{T+2-t}, \dots, y_{T-t+s})$  for  $2 \leq s \leq t$

$\nearrow$   
no such intermediate quantity appears twice

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

\* 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  $G\sqrt{(2\ln N)/T}$

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  $(\psi_t^{\text{EG}}(y_{T+2-t}, \dots, y_T), y_{T+1}) = (\tilde{\psi}_t, y_{T+1})$  have the same distribution, so that

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

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

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

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

same  $\varepsilon_T$  as for the iid case.