

## Sequential setting #1 (meta-statistical)

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

- experts output forecasts  $f_{jt}$ ,  $j \in \{1, \dots, N\}$
- statistician aggregates the forecasts as  $\hat{y}_t = \sum_j p_{jt} f_{jt}$   
where  $p_t = (p_{1t}, \dots, p_{Nt})$  is a convex weight vector
- true observation  $y_t$  is revealed

\* Assessment of performance: loss function  $l(\hat{y}_t, y_t)$

$$\hookrightarrow \text{cumulative loss } \sum_{t=1}^T l(\hat{y}_t, y_t)$$

No stochastic model  $\rightarrow$  relative-performance criterion = perform almost as well as the best constant convex combination

$$\text{ie, control } R_T = \sum_{t=1}^T l(\hat{y}_t, y_t) - \inf_{q \in X} \sum_{t=1}^T l\left(\sum_{j=1}^N q_j f_{jt}, y_t\right)$$

↑  
the simplex  
of all convex weight vectors

↑  
independent of  $t$

Ex: square loss  $l(\hat{y}_t, y_t) = (\hat{y}_t - y_t)^2$

absolute loss  $l(\hat{y}_t, y_t) = |\hat{y}_t - y_t|$

absolute percentage of error  $l(\hat{y}_t, y_t) = |\hat{y}_t - y_t| / |y_t|$

$\hookrightarrow$  all these loss functions are convex in  $\hat{y}_t$ , a fact that we will need later on.

\*  $R_T$  is called the regret

We have some bias/variance decomposition:

$$\sum_{t=1}^T l(\hat{y}_t, y_t) = \inf_q \sum_{t=1}^T l\left(\sum_j q_j f_{jt}, y_t\right) + R_T$$

↑  
cumulative loss

We will focus on the regret in these lectures

approximation error ( $\geq$  how good the experts are)  
↑  
the regret corresponds to a sequential situation error

## Sequential setting #2

(sequential optimization)

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

- mathematician picks  $p_t \in X$
- true loss function  $l_t: X \rightarrow \mathbb{R}$  is revealed

\* Assessment of performance

$$R_T = \sum_{t=1}^T l_t(p_t) - \inf_{q \in X} \sum_{t=1}^T l_t(q)$$

\* Link:  $l_t: p \mapsto l(\sum p_j f_j, y_t)$

The only (small) difference is that the  $f_j$  are only available, in this version, AFTER  $p_t$  is selected.

\* Special case of interest to start with:

LINEAR loss functions

$$l_t(p) = \sum_{j=1}^N p_j l_{jt} \quad \text{where } l_{jt} \stackrel{\text{not.}}{=} l_t(s_j)$$

$l_t$  is given by a vector  $(l_{1t}, \dots, l_{Nt})$ ,  
choosing  $l_t$  is choosing such a vector.

↑  
weight  
on  $j$

Setting #2 /  
Linear losses

→ Exponentially weighted average predictor [EWA]

\* Algorithm :

$$p_1 = (\frac{1}{N}, \dots, \frac{1}{N})$$

$$t \geq 2, \quad p_{jt} = \frac{\exp(-\eta \sum_{s=1}^{t-1} l_{js})}{\sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} l_{ks})}$$

where  $\eta > 0$  is a learning rate.

formula  
abs  
for  
since  
 $\sum_{s=1}^{t-1} l_{ks} = 0$   
by convention

Note : Not all weight on the best experts (it wouldn't work!) but  
most of the weight.

\* Bound : For all strategies picking the linear loss functions  $(l_{1t}, \dots, l_{Nt}) \in [m, M]^N$

$$\begin{aligned} R_T &= \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} - \min_{q \in \mathcal{X}} \sum_{t=1}^T \sum_{j=1}^N q_j l_{jt} \\ &= \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} - \min_{k=1 \dots N} \sum_{t=1}^T l_{kt} \leq \frac{\ln N}{\eta} + \eta \frac{(M-m)^2}{8} \end{aligned}$$

\* PROOF : based on Hoeffding's lemma :

- for random variables  $X \in [m, M]$ ,

$$\forall \eta \in \mathbb{R}, \quad \ln \mathbb{E}[e^{\eta X}] \leq \eta \mathbb{E}[X] + \frac{\eta^2 (M-m)^2}{8}$$

- thus for all convex weight vectors  $q \in \mathcal{X}$ , for all  $l_{jt} \in [m, M]$

$$\forall \eta \in \mathbb{R}, \quad \ln \sum_j q_j e^{-\eta l_{jt}} \leq -\eta \sum_j q_j l_{jt} + \frac{\eta^2 (M-m)^2}{8}$$

In particular,

$$\sum_j p_{jt} l_{jt} \leq -\frac{1}{\eta} \ln \sum_j p_{jt} e^{-\eta l_{jt}} + \frac{\eta^2 (M-m)^2}{8}$$

$$= \ln \overbrace{\sum_j}^t \exp(-\eta \sum_{s=1}^t l_{js})$$

$$\uparrow \sum_k \exp(-\eta \sum_{s=1}^{t-1} l_{ks})$$

by definition  
of the algorithm

Summing over  $t$  (telescoping sum in the right-hand side):

$$\sum_{t=1}^T \sum_{j=1}^N p_j^t l_j^t \leq -\frac{1}{\eta} \ln \frac{\sum_{j=1}^N \exp(-\eta \sum_{t=1}^T l_j^t)}{N} + \frac{\eta}{8} (M-m)^2 T$$

Conclude using  $\ln \sum_j e^{-\eta \sum_{t=1}^T l_j^t} \geq \max_j \ln e^{-\eta l_j^m} = \max_j \ln e^{-\eta l_j^m} = \max_j \{-\eta l_j^m\}$

and (linearity):  $\min_{j \in [1..N]} \sum_{t=1}^T l_j^t = \min_{q \in \mathcal{X}} \sum_{t=1}^T \sum_k q_k l_{k,t}$ .  $\square$

Reminder :

PROOF OF HOEFFDING'S LEMMA:

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

$$\Psi'(\eta) = \frac{\mathbb{E}[X e^{\eta X}]}{\mathbb{E}[e^{\eta X}]} \quad (\text{differentiability: f. } X \text{ bounded locally})$$

$$(\Psi \text{ twice differentiable: similar reasons}) \quad \Psi''(\eta) = \frac{\mathbb{E}[X^2 e^{\eta X}] \mathbb{E}[e^{\eta X}] - (\mathbb{E}[X e^{\eta X}])^2}{(\mathbb{E}[e^{\eta X}])^2} = \text{Var}_{\mathbb{Q}}(X)$$

under the probability  $\mathbb{Q}$  defined by

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

$$X \in [m, M] : \text{Var}_{\mathbb{Q}}(X) = \inf_{\mu \in \mathbb{R}} \mathbb{E}_{\mathbb{Q}}[(X-\mu)^2] \leq \mathbb{E}_{\mathbb{Q}}[(X - \frac{M+m}{2})^2] \leq \frac{(M-m)^2}{4}$$

Taylor:  $\exists x \text{ s.t. } \Psi(\eta) = \overset{\eta=0}{\Psi(0)} + \eta \Psi'(0) + \frac{\eta^2}{2} \overset{\eta=x}{\Psi''(x)}$   
 $(\text{f. } \Psi \text{ is actually } C^2)$

$$\text{i.e. } \ln \mathbb{E}[e^{\eta X}] \leq \eta \mathbb{E}[X] + \frac{\eta^2}{2} (M-m)^2 \quad \square$$

\* How good is the obtained bound?

$$\min_{\eta > 0} \frac{\ln N}{\eta} + \frac{\eta}{8} (M-m)^2 T = (M-m) \sqrt{\frac{T \ln N}{2}} \quad \text{for } \eta^* = \frac{1}{(M-m)} \sqrt{\frac{8 \ln N}{T}}$$

↑ nice because  $O(T)$  and weak dependence in  $N$

↑ not nice when  $T, m, M$  unknown

Why bother? Why not put all mass on the best  $j$  in EWA?  
 (exercise) #1 Why only compare the performance to the best global expert?  
 How good is EWA?

1) Consider the strategy

$$p_t = (\gamma_N \dots \gamma_1)$$

$$p_t \leftrightarrow \begin{cases} \text{equal mass on } k \text{ s.t. } \sum_{s=t+1}^T l_{ks} = \min_j \sum_{s=t+1}^T l_{js} \\ \text{zero mass on } k \text{ s.t. } \sum_{s=t+1}^T l_{ks} > \min_j \sum_{s=t+1}^T l_{js} \end{cases}$$

called : « follows the leader(s) » (FTL)

|| Show that it fails: there exist sequences  $(l_{it}, \dots, l_{nt}) \in [0, 1]^N$   
 such that  $\exists S > 0$ ,  $R_T = \sum_{t=1}^T \sum_{j=1}^N p_{jt} l_{jt} - \min_{k=1-N} \sum_{t=1}^T l_{kt} \geq S$ .

EWA appears as a smoothed version of FTL!

2) Can we be more ambitious? People often hope to get close to  
 $\sum_{t=1}^T \min_{k=1-N} l_{kt}$  in statistics when they resort to aggregation.

|| Show that no strategy can be such that for all sequences  
 $(l_{it}, \dots, l_{nt}) \in [0, 1]^N$ ,  $R_T = \sum_{t,j} p_{jt} l_{jt} - \sum_t \min_k l_{kt} = o(T)$ .

3) Show that for EWA, the regret is never negative:

||  $\forall (l_{it}, \dots, l_{nt}) \in [m, M]^N$ ,  $R_T = \sum_{j,t} p_{jt} l_{jt} - \min_{k=1-N} \sum_t l_{kt} \geq 0$

REMINDER:

When does a limit in distribution entail a limit for the expectations?

(1) In general,  $Y_T \xrightarrow{d} Y$  does not imply  $E[f(Y_T)] \rightarrow E[f(Y)]$   
if  $f$  is just assumed continuous.

Ex:  $f = \text{identity}$ ,  $Y_T = T^2 \mathbb{1}_{\{U_T \in [a^{1/T}]^c\}}$  where  $U_1, \dots, U_T, \dots$  iid  $\sim U_{[0,1]}$   
We have  $\int Y_T \frac{dP}{dP} \rightarrow 0$   
 $E[Y_T] = T \rightarrow +\infty$

(2) Definition of  $\xrightarrow{d}$  is for continuous and bounded functions:

Let  $(Y_T)_{T \geq 1}$  be a sequence of  $\mathbb{R}^N$ -valued random variables

$Y_T \xrightarrow{d} Y$  if by definition, for all continuous and bounded functions  $f: \mathbb{R}^N \rightarrow \mathbb{R}$ , we have

$$E[f(Y_T)] \rightarrow E[f(Y)].$$

(3) Uniform asymptotic integrability [uai]: is a sufficient condition for

Def:  $(Z_T)_{T \geq 1}$  is uai if

$$\lim_{T \rightarrow +\infty} \limsup_{\varepsilon \rightarrow 0} E[\|Z_T\| \mathbb{1}_{\{\|Z_T\| > \varepsilon\}}] = 0$$

Lemma: If  $(Y_T)_{T \geq 1}$  is a sequence of random vectors in  $\mathbb{R}^N$  and  $f: \mathbb{R}^N \rightarrow \mathbb{R}$  is a continuous function such that

- $Y_T \xrightarrow{d} Y$

- $(f(Y_T))_{T \geq 1}$  is uai,

then: a)  $f(Y) \in L^1$  and  $f(Y_T) \in L^1$  for all  $T$  sufficiently large

b)  $E[f(Y_T)] \rightarrow E[f(Y)]$

Proof: For a): by uai,  $\exists L_0 | \limsup_{T \rightarrow +\infty} E[|f(Y_T)| \mathbb{1}_{\{|f(Y_T)| > L_0\}}] < +\infty$   
in particular,  $\limsup_{T \rightarrow +\infty} E[|f(Y_T)|] < +\infty$ ,

thus,  $(E[|f(Y_T)|])_{T \geq T_0}$  is bounded for some  $T_0$

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

$$f(Y_T) \rightarrow f(Y) \text{ as } T \rightarrow \infty$$

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

For b): We denote by  $\Psi_L : x \in \mathbb{R} \mapsto \max\{-L, \min\{x, L\}\}$  the clipping operator

$$\begin{aligned} |E[f(Y)] - E[f(Y_T)]| &\leq |E[f(Y)] - E[\Psi_L(f(Y))]| \\ &+ |E[\Psi_L(f(Y))] - E[\Psi_L(f(Y_T))]| \\ &+ |E[\Psi_L(f(Y_T))] - E[f(Y_T)]| \end{aligned}$$

Thus,

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

as  $\Psi_L \circ f$  is continuous and bounded

Taking  $\limsup_{L \rightarrow \infty}$  in both sides (well, the left-hand side is independent of  $L$ ) and using that  $\limsup_{T \rightarrow \infty} a_T + b_T \leq \limsup_{T \rightarrow \infty} a_T + \limsup_{T \rightarrow \infty} b_T$  for non-negative sequences  $(a_T)$  and  $(b_T)$ :

$$\begin{aligned} \limsup_{T \rightarrow \infty} |E[f(Y)] - E[f(Y_T)]| &\leq \limsup_{L \rightarrow \infty} |E[f(Y)] - E[\Psi_L(f(Y))]| \\ &+ \limsup_{L \rightarrow \infty} \limsup_{T \rightarrow \infty} |E[\Psi_L(f(Y_T))] - E[\Psi_L(f(Y_T))]| \\ &= 0 \text{ by dominated convergence} \quad \text{to be handled by the main property} \end{aligned}$$

By definition of  $\Psi_L$ ,

$$\begin{aligned} \mathbb{E}[f(y_T)] &= \mathbb{E}[\Psi_L(f(y_T))] \\ &= \mathbb{E}\left[(f(y_T) + L) \mathbb{1}_{\{f(y_T) < -L\}}\right] + \mathbb{E}\left[(f(y_T) - L) \mathbb{1}_{\{f(y_T) > L\}}\right] \end{aligned}$$

In particular,

$$\begin{aligned} &|\mathbb{E}[f(y_T)] - \mathbb{E}[\Psi_L(f(y_T))]| \\ &\leq \mathbb{E}\left[\underbrace{|f(y_T) + L|}_{\leq |f(y_T)| \text{ on this event}} \mathbb{1}_{\{f(y_T) < -L\}}\right] + \mathbb{E}\left[\underbrace{|f(y_T) - L|}_{\leq |f(y_T)| \text{ on this event}} \mathbb{1}_{\{f(y_T) > L\}}\right] \\ &\leq \mathbb{E}[|f(y_T)| \mathbb{1}_{\{|f(y_T)| > L\}}] \end{aligned}$$

whose  $\limsup_{L \rightarrow \infty} \limsup_{T \rightarrow \infty} \dots$  equals 0 by the very definition of uai for  $(f(y_T))_{T \geq 1}$ .

This concludes the proof.

(4) A useful sufficient condition for uai:

Lemma: If  $(z_T)_{T \geq 1}$  is bounded in  $\mathbb{L}^p$  for some  $p > 1$ ,  
i.e.  $\sup_{T \geq 1} \mathbb{E}[|z_T|^p] = B < +\infty$ .

Then  $(z_T)_{T \geq 1}$  is uai.

Proof:  $\varphi(x) = x^p$ : as  $p > 1$ ,  $\varphi(x) \geq x$

thus  $\forall M > 0$ ,  $\exists L \mid \forall x \geq L$ ,  $\varphi(x) \geq Mx$

$$\begin{aligned} \text{Taking } M = Bg: \exists L \mid \forall T \geq 1, B > \mathbb{E}[|z_T|^p] \geq \mathbb{E}[|z_T|^p \mathbb{1}_{\{|z_T| > Lg\}}] \\ \geq Bg \mathbb{E}[|z_T| \mathbb{1}_{\{|z_T| > Lg\}}] \end{aligned}$$

Thus, for this  $L$ :

$$\sup_{T \geq 1} \mathbb{E}[|z_T| \mathbb{1}_{\{|z_T| > Lg\}}] \leq \frac{1}{g}$$

and the uai property follows by letting  $g \rightarrow +\infty$ .  
(even something stronger than that follows)

OPTIMALITY OF THE

$$\sqrt{\frac{T}{2} \ln N}$$

BOUND:

Asymptotic lower bound

- In the case of linear losses

- By homogeneity we may assume  $\text{Pct } l_{tj} \in [0,1]$  (ie,  $m=0$  and  $M=1$ ) for all  $t$  and  $j$

Theorem:

$$\liminf_{N \rightarrow +\infty} \liminf_{T \rightarrow +\infty} \inf_{\text{algorithms}} \sup_{l_{tj} \in [0,1]} \frac{\sum_{t,j} p_{jt} l_{jt} - \min_k \sum_t l_{kt}}{\sqrt{\frac{T}{2} \ln N}} \geq 1$$

↑  
here stronger to have an infimum than to prove  
algo,  $\liminf \liminf \sup \dots \geq 1$

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

Proof:

We lower bound the  $\sup_{l_{tj} \in [0,1]}$  by an  $E[\dots]$

with losses  $l_{jt}$  drawn iid  $\sim \text{Ber}(\frac{1}{2})$

Using Pct for any algorithm:

$$\forall t \geq 1 \quad E\left[ \sum_j p_{jt} l_{jt} \right] = \sum_j E[p_{jt}] E[l_{jt}] = \frac{1}{2} E\left[ \sum_j p_{jt} \right] = \frac{1}{2}$$

$p_{jt}$  is measurable w.r.t past losses  $L_{ks}$   $s \leq t-1$  and is thus independent of the  $L_{jt}$

↳ Not so surprising: means Pct we cannot predict efficiently fair coin tosses

\* But \* by a crowd effect (due to the central limit theorem)  
we will show that

$$E\left[ \min_k \sum_{t=1}^T L_{kt} \right] = T_2 - \text{something of order } \frac{1}{\sqrt{T \ln N}}$$

even though we "individually" have

$$\forall k, E\left[ \sum_{t=1}^T L_{kt} \right] = T_2$$

More formally:

$$\inf_{\text{algo}} \sup_{l_{tj} \in [0,1]} \frac{\sum_{t,j} p_{jt} l_{jt} - \min_k \sum_t l_{kt}}{\sqrt{\frac{T}{2} \ln N}}$$

$$= \frac{1}{\sqrt{2 \ln N}} E\left[ \max_{k=1 \dots N} \frac{\sum_{t=1}^T (T_2 - L_{kt})}{\frac{1}{2} \sqrt{T}} \right]$$

$$\frac{E\left[ \sum_{t,j} p_{jt} l_{jt} \right] - E\left[ \min_k \sum_{t=1}^T L_{kt} \right]}{\sqrt{(T_2) \ln N}}$$

We denote

$$Z_{kT} = \frac{\sum_{t=1}^T (\frac{1}{2} - L_{kt})}{\frac{1}{2}\sqrt{T}}$$

← sum of iid random var  
←  $\sigma\sqrt{T}$

and note that  $Z_{kT}$  is a central-limit-theorem statistic:  $\frac{1}{2} - L_{kt}$  are iid, centered with  $\text{Var}(\frac{1}{2} - L_{kt}) = 1/4$

For each  $k$ ,  $Z_{kT} \xrightarrow{d} \mathcal{N}(0, 1)$

By independence,

$$Z_T = \begin{pmatrix} Z_{1T} \\ \vdots \\ Z_{NT} \end{pmatrix} \xrightarrow{d} Z \sim \mathcal{N}(0, I_N)$$

Take  $f(x_1, \dots, x_N) = \max_{k \leq N} x_{k,j}$  we have shown so far:

$$\inf_{\text{algo}} \sup_{f \in \mathcal{F}[a_1]} \frac{\sum_{t=1}^T p_t f_{jt} - \min_{k \leq T} \sum_{t=1}^T l_{kt}}{\sqrt{\sum_{t=1}^T \ln N}} \geq \frac{1}{\sqrt{2 \ln N}} E[f(Z_T)]$$

The  $f(Z_T)$  are bounded by  $N$  in  $L^2$ -norm:

$$E[f(Z_T)^2] = E[\max_{k \leq N} Z_{kT}^2] \leq \underbrace{\sum_{k=1}^N E[Z_{kj}^2]}_{\text{each } = 1} = N$$

cf.  $Z_{kT}$  central-limit-theorem statistic

Thus they form a uoi sequence and we have

$$\lim_{T \rightarrow +\infty} \frac{1}{\sqrt{2 \ln N}} E[f(Z_T)] = \frac{1}{\sqrt{2 \ln N}} E[\max_{k \leq N} Z_k]$$

So far:

$$\liminf_{T \rightarrow +\infty} \inf_{\text{algo}} \sup_{f \in \mathcal{F}[a_1]} \frac{\sum_{t=1}^T p_t f_{jt} - \min_{k \leq T} \sum_{t=1}^T l_{kt}}{\sqrt{(T/2) \ln N}}$$

$$\geq \frac{1}{\sqrt{2 \ln N}} E[\max_{k \leq N} Z_k]$$

                

The proof is concluded by:

↑ details on the next pages

$$\liminf_{N \rightarrow +\infty} \dots \geq 1$$

(This is even a true limit)

## Asymptotics for the expectation of the maximum of standard Gaussian random variables

We prove that for  $G_1, G_2, \dots, G_N$  iid  $\sim \mathcal{N}(0, 1)$ ,

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

(1) Easy part (useless for the main result proved above):

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

as a special case of the following lemma.

Lemma: If  $V_1, V_2, \dots, V_N$  are such that  $\mathbb{E}[e^{dV_j}] \leq e^{d^2 \sigma^2 / 2}$  for all  $j$ , then  $\mathbb{E}[\max_{k \leq N} V_k] \leq \sqrt{2\sigma^2 \ln N}$

Notes: • Valid even without independence (independence will be needed for the lower bound)

• If  $G_k \sim \mathcal{N}(0, \sigma^2)$  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 by Hoeffding's lemma

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

• We call sub-Gaussian the random variables  $V_k$  satisfying the assumptions of the lemma.

Proof: Proof technique called "Pisier's argument":

$$\mathbb{E}[\max_{k \leq N} V_k] \leq \frac{1}{d} \ln \mathbb{E}[\exp(d \max_{k \leq N} V_k)] = \frac{1}{d} \ln \mathbb{E}[\max_{k \leq N} e^{dV_k}]$$

$$\stackrel{\text{↑ Jensen's inequality}}{\leq} \frac{1}{d} \ln \left( \sum_{k \leq N} \mathbb{E}[e^{dV_k}] \right)$$

$$\stackrel{\text{↑ max of non-neg quantities}}{\leq} \frac{1}{d} \ln \left( N e^{d^2 \sigma^2 / 2} \right) = \frac{\ln N}{d} + \frac{d \sigma^2}{2}$$

$$\stackrel{\text{↑ by the sub-Gaussian assumption}}{=} \sqrt{2\sigma^2 \ln N}$$

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

(2) “Difficult” part (and the one that we need!):

Credits: Pascal Nassart +  
Mathieu Even,  
M2 student following  
this course in 2019-20

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

$$M_N \stackrel{\text{def.}}{=} \left( \max_{k \leq N} G_k / \sqrt{2 \ln N} \right)^+ \quad \leftarrow \text{non-negative part}$$

$$\begin{aligned} \text{Given that } 0 &\geq \mathbb{E} \left[ \left( \max_{k \leq N} G_k \right) \mathbf{1}_{\{\max_{k \leq N} G_k \leq 0\}} \right] \\ &\geq \mathbb{E} \left[ G_1 \mathbf{1}_{\{\max_{k \leq N} G_k \leq 0\}} \right] \geq \mathbb{E} \left[ G_1 \mathbf{1}_{\{G_1 \leq 0\}} \right] \\ &\geq -\sqrt{\mathbb{E}[G_1^2] \mathbb{P}[G_1 \leq 0]} = -\sqrt{2} \end{aligned}$$

(Cauchy-Schwarz)

we have

$$\liminf_{N \rightarrow \infty} \frac{\mathbb{E}[\max_{k \leq N} G_k]}{\sqrt{2 \ln N}} = \liminf_{N \rightarrow \infty} \mathbb{E}[M_N]$$

By Fatou's lemma, it thus suffices to show that  $\liminf_{N \rightarrow \infty} M_N \geq 1$  as.

For  $\varepsilon \in (0, 1)$ :

$$\begin{aligned} \mathbb{P}(M_N \leq \sqrt{1-\varepsilon}) &= \mathbb{P}\left(\forall k \leq N, G_k \leq \sqrt{2(1-\varepsilon)\ln N}\right) \\ &= (\Phi(\sqrt{2(1-\varepsilon)\ln N}))^N \underset{\text{true}}{\leq} \exp(-N(1-\Phi(\dots))) \end{aligned}$$

where  $\Phi$  is the cdf of the  $\mathcal{N}(0, 1)$  distribution.

$$\text{As } 1 - \Phi(x) \underset{x \rightarrow \infty}{\sim} e^{-x^2/2} / x \sqrt{2\pi} \quad (\text{proof by two integration by parts, see below})$$

we have

$$\begin{aligned} N(1 - \Phi(\sqrt{2(1-\varepsilon)\ln N})) &\sim N \exp(-((1-\varepsilon)\ln N) / 2\sqrt{\pi(1-\varepsilon)\ln N}) \\ &= N^\varepsilon / 2\sqrt{\pi(1-\varepsilon)\ln N} \end{aligned}$$

Thus  $\sum_N \exp(-N(1 - \Phi(\dots)))$  is convergent

therefore,

$\sum_N \mathbb{P}(M_N \leq \sqrt{1-\varepsilon})$  is also convergent

By the Borel-Cantelli lemma:

$$\{M_N \leq \sqrt{1-\varepsilon}\}$$

a.s. takes place only finitely many times

i.e.  $\liminf_{N \rightarrow \infty} M_N \geq \sqrt{1-\varepsilon}$  a.s.

Since this is true for all  $\varepsilon > 0$ , we get the desired  $\liminf_{N \rightarrow \infty} M_N \geq 1$  a.s. (with  $\varepsilon_n = \eta_n \rightarrow 0$ )

(3) Proof of  $1 - \Phi(x) \underset{x \rightarrow +\infty}{\sim} e^{-x^2/2} / \sqrt{2\pi}$

FIRST IBP:  
(integration by parts)

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

Second IBP:

$$\begin{aligned} 0 &\leq \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} \frac{e^{-t^2/2}}{t^2} dt = \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} t e^{-t^2/2} \frac{1}{t^3} dt \\ &= \underbrace{\left[ -\frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t^3} \right]_x^{+\infty}}_{\leq \frac{1}{\sqrt{2\pi}} \frac{e^{-x^2/2}}{x^3}} - \underbrace{\int_x^{+\infty} dt}_{> 0} \end{aligned}$$

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

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  $\Pi_0$  be a probability such that the losses  $l_{kt}$  are iid  $\sim \text{Ber}(\frac{1}{2})$

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

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

$$(1) \text{ Show that } \sup_{\Pi_T \in \{\Pi_0, \dots, \Pi_N\}} \left\{ \sum_{j=1}^N p_{jt} l_{jt} - \min_k \sum_{t=1}^T l_{kt} \right\} \geq T \cdot \max_{k=1, \dots, N} \left\{ 1 - \frac{1}{T} \sum_{t=1}^T E_k [p_{kt}] \right\}$$

REMINDER or QUICK INTRODUCTION TO A NEW OBJECT:

Recall that the Kullback-Leibler divergence between two probability distributions

$\mu$  and  $\nu$  equals:

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

KL has the following properties:

- $\text{KL}(\nu \otimes \nu', \mu \otimes \mu') = \text{KL}(\nu, \mu) + \text{KL}(\nu', \mu')$  ↗ by Fubini-Tonelli
- $\forall$  random variable  $X$ ,  $\text{KL}(\nu^X, \mu^X) \leq \text{KL}(\nu, \mu)$  ↗ data-processing inequality  
(we take it for granted and will prove it later)

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

$$= p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} \quad \text{Re KL-}$$

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

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

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

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:

- $KL(pq) \leq K$  / to that end, first prove that  $KL$  is jointly convex (it is a consequence of the data processing inequality, explain why)
- $KL(pq)/\ln N \geq a$  for  $a \geq 2e/(2e+1)$

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

Show

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

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

with  $L = (L_{kt})_{\substack{k \in [N] \\ t \in [T]}}$

(4) Show that  $K' \leq 5T\Sigma^2$  for  $\Sigma \leq \gamma_0$ .

(5) Conclude to the following bound:

Theorem: For all strategies, for all  $N \geq 2$  and  $T \geq 17 \ln N$ ,

$$\sup R_T \geq 0.06 \sqrt{T \ln N}$$

Application of the EWNA forecaster / Sion's lemma.

Statement: Let  $X, Y$  two convex sets,  $f: X \times Y \rightarrow [0, M]$

a function s.t.

$\forall x \in X$ ,  $f(x, \cdot)$  is concave

$\forall y \in Y$ ,  $f(\cdot, y)$  is convex

then (under additional regularity assumptions):

$$\inf_{x \in X} \sup_{y \in Y} f(x, y) = \sup_{y \in Y} \inf_{x \in X} f(x, y).$$

Proof:

1)  $\geq$  always holds :  $\forall x, y$ ,  $f(x, y) \geq \inf_{x' \in X} f(x', y)$

taking  $\sup_{y \in Y}$  in both sides:  $\forall y \sup_{y \in Y} f(x, y) \geq \sup_{y \in Y} \inf_{x \in X} f(x, y)$

(Get  $\inf_x \sup_y f \geq \sup_x \inf_y f$  by taking the  $\inf_{x \in X}$  (the right-hand side is a constant independent of  $x$ )).

2) A (fictitious) statistician and a (fictitious) opponent play as follows :

First, the statistician sets  $N \geq 2$  and  $x^{(1)}, \dots, x^{(N)}$  in  $X$ , as well as  $T \geq 1$ .

Then, at each round, they simultaneously pick

$$x_t = \sum_{j=1}^N p_j^{(t)} x^{(j)} \in X \quad \text{and} \quad y_t \in Y$$

How?  $p_j^{(t)} = \exp(-\eta \sum_{s=1}^{t-1} f(x^{(s)}, y_s)) / \sum_{k=1}^N \exp(-\eta \sum_{s=1}^{t-1} f(x^{(k)}, y_s))$

with  
 $\eta = \frac{1}{M} \sqrt{\frac{\ln N}{T}}$

and (since  $p_j^{(t)}$  only depends on the past, the opponent can compute it and pick :)

$y_t$  st.  
(by definition of sup)

$$f(x_t, y_t) \geq \sup_{y \in Y} f(x_t, y)$$

$$-\frac{1}{T}$$

By definition of the exponentially weighted average strategy with a well-chosen learning rate  $\eta \geq 0$

$$\sum_{t=1}^T \underbrace{\sum_{j=1}^n p_j f(x^{(j)}, y_t)}_{\text{corresponds to } \min_{k \in [0, M]} f(x^{(k)}, y_t)} - \min_{k \in [0, N]} \sum_{t=1}^T f(x^{(k)}, y_t) \leq \underbrace{\frac{\ln N}{\eta}}_{\text{given our choice for } \eta} + \eta \frac{M^2}{8} T = M \sqrt{\frac{T}{2} \ln N}$$

From the convexity of  $f(\cdot, y_t)$ , we finally get:

$$\sum_{t=1}^T f\left(\underbrace{\frac{1}{j} \sum_{j=1}^n x^{(j)}}_{= \bar{x}_t \text{ def.}}, y_t\right) - \min_{k \in [0, N]} \sum_{t=1}^T f(x^{(k)}, y_t) \leq M \sqrt{\frac{T}{2} \ln N} \quad (*)$$

3) Now,  $\inf_x \sup_y f(x, y) \leq \sup_y f(\bar{x}, y)$

$$\leq \sup_y \frac{1}{T} \sum_{t=1}^T f(x_t, y) \stackrel{\sup \frac{1}{T} \sum}{\leq} \frac{1}{T} \sum_{t=1}^T \sup_y f(x_t, y) \stackrel{\text{def of } y_t}{\leq} \frac{1}{T} \sum_{t=1}^T f(x_t, y_t) + \frac{1}{T}$$

while by (\*):

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T f(x_t, y_t) \\ & \leq \underbrace{M \sqrt{\frac{\ln N}{2T}}}_{= O(\sqrt{\chi_T})} + \min_{k \in [0, N]} \frac{1}{T} \sum_{t=1}^T f(x^{(k)}, y_t) \stackrel{\text{by } (*)}{\leq} \\ & \leq O(\sqrt{\chi_T}) + \min_k \underbrace{f(x^{(k)}, \bar{y})}_{\leq \sup_{y \in Y} \min_k f(x^{(k)}, y)} \stackrel{\text{by concavity of } f(x^{(k)}, \cdot), \text{ where } \bar{y} = \frac{1}{T} \sum_{t=1}^T y_t}{=} \end{aligned}$$

4) In sections (2)-(3)  $T, N$  and  $x^{(1)}, \dots, x^{(n)}$  were fixed, but we can play with them! We proved

$$\inf_x \sup_y f(x, y) \leq \frac{1}{T} + O(\sqrt{\chi_T}) + \sup_y \min_k f(x^{(k)}, y)$$

Letting  $T \rightarrow +\infty$ :

$$\inf_x \sup_y f(x, y) \leq \sup_y \min_k f(x^{(k)}, y)$$

This holds for all  $N$  and all  $x^{(1)} \dots x^{(N)}$  in  $X$ :

$$\inf_x \sup_y f(x, y) \leq \inf_{N \geq 1} \inf_{\{x^{(1)}, \dots, x^{(N)}\} \subset X} \sup_{y \in Y} \min_{k=1 \dots N} f(x^{(k)}, y)$$

$\geq \sup_y \inf_x f(x, y)$  but maybe  $\leq \sup_y \inf_x f(x, y)$  so that we have an equality?  
 We will now state and use regularity / topological assumptions.

Assume •  $X, Y$  are metric spaces, with distances  $d_X$  and  $d_Y$

•  $f: X \times Y \rightarrow [0, M]$  uniformly continuous:

(In particular true when  $X$  is included in a compact set)

$$\forall \varepsilon > 0, \exists \delta > 0 \mid d_X(x_1, x_2) + d_Y(y_1, y_2) \leq \delta$$

$$\Rightarrow |f(x_1, y_1) - f(x_2, y_2)| \leq \varepsilon$$

• Finite covering property for  $X$ :

$$\forall \delta > 0, \exists N \text{ and } x^{(1)} \dots x^{(N)} \text{ s.t. } X \subset \bigcup_{j=1}^N B(x^{(j)}, \delta)$$

Given  $\varepsilon > 0$  and the associated  $\delta > 0$ :

$$\forall x, y, f(x, y) \geq \min_{j=1 \dots N} f(x^{(j)}, y) - \varepsilon \quad \text{by uniform continuity}$$

Taking  $\inf_x$  then  $\sup_y$ :

$$\sup_y \inf_x f(x, y) \geq \sup_y \min_j f(x^{(j)}, y) - \varepsilon$$

$$\geq \inf_{\{x^{(j)}\}} \sup_y \min_j f(x^{(j)}, y)$$

and we let  $\varepsilon \downarrow 0$   $- \varepsilon$

Exercise #2

Bonus points rewarded!

(Perhaps you can find even better - weaker - assumptions? If so, let me know!)

↑ while still having a smooth and easy-to-read proof...)