

Part 3: Fano-type inequalities.

Simplest framework: measurable space (Ω, \mathcal{F})
 finitely many probability distributions $\mathbb{P}_i, \mathbb{Q}_i$ $i \in \{1, \dots, N\}$
 events $A_i \in \mathcal{F}$

Aim: bound $\frac{1}{N} \sum_{i=1}^N \mathbb{P}_i(A_i)$ by $\frac{1}{N} \sum_{i=1}^N \mathbb{Q}_i(A_i) + \dots$

Most general framework: underlying measurable space (Ω, \mathcal{F})
 parameter space (Θ, \mathcal{G}) with prior distribution ν
 families of probability distributions $\mathbb{P}_\theta, \mathbb{Q}_\theta$ $\theta \in \Theta$
 family of $[a, b]$ -valued random variables Z_θ $\theta \in \Theta$
 (+ measurability assumptions: $\theta \mapsto E_{\mathbb{P}_\theta}[Z_\theta]$ and $\theta \mapsto E_{\mathbb{Q}_\theta}[Z_\theta]$ are \mathcal{G} -measurable)

Aim: bound $\int_{\Theta} E_{\mathbb{P}_\theta}[Z_\theta] d\nu(\theta)$ by $\int_{\Theta} E_{\mathbb{Q}_\theta}[Z_\theta] d\nu(\theta) + \dots$

- Two steps:
1. Reduction to Bernoulli distributions
 2. Lower bound the KL function

Reminder of a short-hand notation: $k(p, q) = KL(\text{Ber}(p), \text{Ber}(q))$

Step 1:

$$k\left(\frac{1}{N} \sum_{i=1}^N \mathbb{P}_i(A_i), \frac{1}{N} \sum_{i=1}^N \mathbb{Q}_i(A_i)\right) \leq \frac{1}{N} \sum_{i=1}^N k(\mathbb{P}_i(A_i), \mathbb{Q}_i(A_i))$$

\swarrow convexity of KL
 \searrow data-processing inequality

$$\leq \frac{1}{N} \sum_{i=1}^N KL(\mathbb{P}_i, \mathbb{Q}_i)$$

$$k\left(\int_{\Theta} E_{\mathbb{P}_\theta}[Z_\theta] d\nu(\theta), \int_{\Theta} E_{\mathbb{Q}_\theta}[Z_\theta] d\nu(\theta)\right) \leq \int_{\Theta} k(E_{\mathbb{P}_\theta}[Z_\theta], E_{\mathbb{Q}_\theta}[Z_\theta]) d\nu(\theta)$$

$$\leq \int_{\Theta} KL(\mathbb{P}_\theta, \mathbb{Q}_\theta) d\nu(\theta)$$

(under the needed measurability assumptions)

In both cases: $K(\bar{p}, \bar{q}) \leq \bar{K}$

Step 2:
$$\begin{cases} K(p, q) \geq p \ln \frac{1}{q} - \ln 2 \\ K(p, q) \geq \left(\ln \frac{1}{q}\right) (p-q)^2 \end{cases} \text{ for all } (p, q) \in [0, 1]^2$$

as we saw last time -

Thus:
$$\begin{cases} p \leq \frac{\ln 2 + K(p, q)}{-\ln(1/q)} \\ p \leq q + \sqrt{\frac{K(p, q)}{\ln(1/q)}} \end{cases} \text{ for all } p \in [0, 1] \text{ and } q \in (0, 1)$$

Conclusions:

E.g.: (1)
$$\frac{1}{N} \sum_{i=1}^N P_i(A_i) \leq \frac{\ln 2 + \frac{1}{N} \sum_{i=1}^N K(P_i, Q_i)}{-\ln\left(\frac{1}{N} \sum_{i=1}^N Q_i(A_i)\right)}$$

(2)
$$\frac{1}{N} \sum_{i=1}^N P_i(A_i) \leq \frac{1}{N} \sum_{i=1}^N Q_i(A_i) + \sqrt{\frac{\sum_{i=1}^N K(P_i, Q_i)}{-\ln\left(\frac{1}{N} \sum_{i=1}^N Q_i(A_i)\right)}}$$

(3)
$$\int_{\Theta} E_{P_{\Theta}}[Z_{\Theta}] dJ(\Theta) \leq \int_{\Theta} E_{Q_{\Theta}}[Z_{\Theta}] dJ(\Theta) + \sqrt{\frac{\int_{\Theta} K(P_{\Theta}, Q_{\Theta}) dJ(\Theta)}{-\ln \int_{\Theta} E_{Q_{\Theta}}[Z_{\Theta}] dJ(\Theta)}}$$

Comments:

If the A_i form a partition of Ω and $Q_i = Q \quad \forall i$:

(1) says:
$$\frac{1}{N} \sum_{i=1}^N P_i(A_i) \leq \frac{\ln 2 + \frac{1}{N} \sum_{i=1}^N K(P_i, Q)}{\ln N}$$

→ non-trivial only if $N \geq 3$

Birge's (2005) version:
$$\min_{i=1, \dots, N} P_i(A_i) \leq \max\left\{0.84, \frac{\frac{1}{N-1} \sum_{i=2}^N K(P_i, P_i)}{\ln N}\right\}$$

Alternative idea: resort to (2)!

which indicates
$$\frac{1}{N} \sum_{i=1}^N P_i(A_i) \leq \frac{1}{N} + \sqrt{\frac{\frac{1}{N} \sum_{i=1}^N K(P_i, Q)}{\ln N}}$$

→ Which Q should be chosen in the case $Q_i = Q \forall i$?

Left-hand side independent of Q

Right-hand side: minimize $\frac{1}{N} \sum_{i=1}^N KL(P_i, Q)$

See Yang and Barron (1999):

We may restrict our attention to Q s.t. $P_i \ll Q \forall i$

then $\bar{P} = \frac{1}{N} \sum_{i=1}^N P_i$ is s.t. $\forall j, P_j \ll \bar{P} \ll Q$

in particular, Q -as: $\frac{dP_j}{dQ} = \frac{dP_j}{d\bar{P}} \frac{d\bar{P}}{dQ}$

& that

$$\begin{aligned} \frac{1}{N} \sum_{j=1}^N KL(P_j, Q) &= \frac{1}{N} \sum_{j=1}^N \int \left(\ln \frac{dP_j}{d\bar{P}} + \ln \frac{d\bar{P}}{dQ} \right) dP_j \\ &= \frac{1}{N} \sum_{j=1}^N KL(P_j, \bar{P}) + KL(\bar{P}, Q) \end{aligned}$$

The minimum over Q is achieved at $Q = \bar{P}$.

Application #1:

Robust sequential learning with sparse losses.

where U_1, U_2, \dots independent

Setting:

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

- The environment picks $(\ell_{1t}, \dots, \ell_{Nt}) \in [0, 1]^N$ while simultaneously, the decision maker picks $I_t \in \{1, \dots, N\}$, possibly at random thanks to $U_t \sim U_{[0,1]}$
- Both players observe $(\ell_{1t}, \dots, \ell_{Nt})$ and I_t

Aim of the decision-maker: minimize, ie, control from above, the regret:

$$R_T = \sum_{t=1}^T \ell_{I_t, t} - \min_{k=1, \dots, N} \sum_{t=1}^T \ell_{kt}$$

Remark:

« Robust sequential learning » because the ℓ_{jt} are not realizations of some well behaved stochastic process, they can be picked by an opponent to annoy the decision maker
See the monograph by Cesa-Bianchi and Lugosi (2006).

Variant / generalization by Kwon and Perchet (2015):

Only s components of $(\ell_{1t}, \dots, \ell_{Nt})$ are non-zero, where $s \in \{1, \dots, N\}$ is fixed and known to the decision maker.

Denote by $H_{t-1} = ((\ell_{js})_{j \in \{1, \dots, N\}}, U_s)_{s \in \{1, \dots, t-1\}}$ the information available to the decision-maker when picking the law of I_t .

Theorem: There exists a strategy for the decision maker [namely, the exponentially weighted average algorithm with a well-chosen time-adaptive learning rate $(\eta_t)_t$] such that: for all $T \geq 1$, for all strategies of the environment leading to losses $\ell_{jt} \in [0, 1]$ such that $(\ell_{1t}, \dots, \ell_{Nt})$ has at most s non-null components,

$$\bar{R}_T = \sum_{t=1}^T \mathbb{E}[\ell_{I_t, t} | \sigma(H_{t-1})] - \min_{k=1, \dots, N} \sum_{t=1}^T \ell_{kt} \leq C \left(\sqrt{sT \frac{\ln N}{N}} + \ln N \right)$$

where C is some numerical constant.

(In Kwon and Perchet, 2015 they use a fixed η tuned as a function of N, s and T and get the result for $C = \sqrt{2}$.)

Consequences:

* The Hoeffding-Azuma inequality ensures that:

$\forall T$, with probability at least $1-\delta$,

$$R_T = \sum_{t=1}^T \ell_{I_t} - \min_{k=1..N} \sum_{t=1}^T \ell_{kt} \leq C \left(\sqrt{ST \frac{\ln N}{N}} + \ln N \right) + \sqrt{\frac{T}{2} \ln \frac{1}{\delta}}$$

* By the tower rule:

E w.r.t decision-maker's internal randomization U_1, U_2, \dots

$$\forall T, \quad \mathbb{E}[R_T] \leq C \left(\sqrt{ST \frac{\ln N}{N}} + \ln N \right).$$

We now turn to the lower bound.

Interestingly it will be proved by taking iid vectors $(\ell_{1t}, \dots, \ell_{Nt})$: it will be a purely stochastic lower bound, while the upper bound is deterministic and robust.

Theorem. Given $s \in [1..N]$, for all strategies of the decision maker, for all $N \geq 2$ and all $T > N \ln N / 16s$, there exists a fixed-in-advance sequence of loss vectors $(\ell_{1t}, \dots, \ell_{Nt}) \in [0, 1]^N$ with each at most s non-zero components, such that

w.r.t the internal randomization

$$\mathbb{E}[R_T] = \sum_{t=1}^T \mathbb{E}[\ell_{I_t}] - \min_{k=1..N} \sum_{t=1}^T \ell_{kt} \geq \frac{1}{32} \sqrt{ST \frac{\ln N}{N}}$$

Credits: Asymptotic proof technique:
Car. Brianchi, Freund, Haussler, Helmbold, Schapire, Warmuth (1997)
Kwon and Perchet (2015) follows it and adapt it because of some correlations

We (Gerchinovitz, Nemrod, Stoltz, 2017?) start as in Kwon and Perchet till asymptotic (central-limit-theorem) arguments are used, at which point we resort to Fano's inequality.

Proof: We fix $\varepsilon \in (0, s/2N)$.

For $i \in \{1, \dots, N\}$, \mathbb{P}_i is such that the random vectors (L_{1t}, \dots, L_{Nt}) are iid as t varies, with common distribution given by:

- two-step randomization
- S_t is a uniform random subset of $\{1, \dots, N\}$ with s elements
 - the conditional law of L_{jt} given S_t is:
 - * δ_0 , Dirac mass at 0 if $j \notin S_t$
 - * $\text{Ber}(1/2)$ if $j \in S_t$ and $j \neq i$
 - * $\text{Ber}(1/2 - \varepsilon \frac{N}{s})$ if $j = i$ and $i \in S_t$
- cf. sparsity constraint

under \mathbb{P}_i , L_{it} is "favored"

Marginal distributions:

$$\begin{cases} \mathbb{P}_i L_{jt} = \text{Ber}(s/2N) & \text{if } j \neq i \\ \mathbb{P}_i L_{it} = \text{Ber}(s/2N - \varepsilon) & \text{if } j = i \end{cases}$$

Thus

$$\begin{aligned} \mathbb{E}_i \left[\mathbb{E}[R_T] \right] &= \mathbb{E}_i \left[\sum_{t=1}^T \mathbb{E}[L_{I_t t}] - \min_{k=1 \dots N} \sum_{t=1}^T L_{kt} \right] \\ &\geq \mathbb{E}_i \otimes \mathbb{E} \left[\sum_{t=1}^T L_{I_t t} \right] - \min_{k=1 \dots N} \underbrace{\mathbb{E}[L_{kt}]}_{= \mathbb{E}[L_{it}] = T \left(\frac{s}{2N} - \varepsilon \right)} \end{aligned}$$

\uparrow w.r.t the decision-maker's internal randomization $U_1, U_2, \dots \text{ iid } \sim \mathcal{U}_{[0,1]}$
 \uparrow w.r.t the L_{jt}

By the tower rule: (applied twice)

$$\begin{aligned} \mathbb{E}_i \otimes \mathbb{E} \left[L_{I_t t} \mid \sigma(H_{t-1}), L_{1t}, \dots, L_{Nt} \right] &= \sum_{j=1 \dots N} L_{jt} \underbrace{\mathbb{P}\{I_t = j \mid \sigma(H_{t-1})\}}_{\text{independent of } \sigma(H_{t-1})} \\ \mathbb{E}_i \otimes \mathbb{E} \left[L_{I_t t} \mid \sigma(H_{t-1}) \right] &= \sum_{j=1 \dots N} \left(\frac{s}{2N} - \varepsilon \mathbb{1}_{\{j=i\}} \right) \times \mathbb{P}\{I_t = j \mid \sigma(H_{t-1})\} \end{aligned}$$

this (we integrate over the L_{jt})

and finally,

$$\mathbb{E}_i \otimes \mathbb{E} [L_{I_t t}] = \frac{s}{2N} - \varepsilon \mathbb{P}_i \otimes \mathbb{P}\{I_t = i\}$$

Summing over $t=1, \dots, T$:

$$\begin{aligned} E_i [E[R_T]] &\geq T \frac{s}{2N} - \varepsilon \sum_{t=1}^T P_i \{ \mathbb{I}_{E=i} \} - (T(\frac{s}{2N} - \varepsilon)) \\ &= T \varepsilon \left(1 - E_i \{ E \left[\frac{N_i(T)}{T} \right] \right) \end{aligned}$$

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

Now, $\underbrace{\sup_{(L_1, \dots, L_T) \in \{0,1\}^N}_{t=1, \dots, T} E[R_T]}_{\text{our aim}} \geq \max_{i \in \{1, \dots, N\}} E_i [E[R_T]] \geq \frac{1}{N} \sum_{i=1}^N E_i [E[R_T]] \geq T \varepsilon \left(1 - \frac{1}{N} \sum_{i=1}^N E_i \{ E \left[\frac{N_i(T)}{T} \right] \} \right)$

The main difference w.r.t. Kwon and Perchet arises now: we replace a central limit theorem argument by an application of Fano's inequality.

\mathcal{Q} such that $\left\{ \begin{array}{l} S_t \text{ is a uniform random subset of } \{1, \dots, N\} \text{ with } s \text{ elements} \\ \text{the conditional law of } L_{jt} \text{ given } S_t \text{ is} \end{array} \right.$

- * $\mathbb{I}_{j \in S_t} \text{ Bern}(1/2)$
- * $\mathbb{I}_{j \notin S_t} \text{ Bern}(1/2)$

same definition as the \mathbb{P}_i , except that no component is favored

Fano's lemma:

$$\frac{1}{N} \sum_{i=1}^N E_i \{ E \left[\frac{N_i(T)}{T} \right] \} \leq \underbrace{\frac{1}{N} \sum_{i=1}^N E_{\mathcal{Q}} \{ E \left[\frac{N_i(T)}{T} \right] \}}_{=1} + \sqrt{\frac{\frac{1}{N} \sum_{i=1}^N \text{KL}(\mathbb{P}_i^T \otimes U^T, \mathcal{Q}^T \otimes U^T)}{-\ln \left(\frac{1}{N} \sum_{i=1}^N E_{\mathcal{Q}} \{ E \left[\frac{N_i(T)}{T} \right] \} \right)}}$$

(we apply Fano not on the underlying probability space, but on the image space of the L_{jt} and U_{jt})

where $\mathbb{P}_i^T = \mathbb{P}_i^{((L_{1t}, \dots, L_{Nt})_{t=1, \dots, T})}$
 $= (\mathbb{P}_i^{(L_{1T}, \dots, L_{NT})})^{\otimes T}$
 and $\mathcal{Q}^T = (\mathcal{Q}^{(L_{1T}, \dots, L_{NT})})^{\otimes T}$
 and $U^T = (U_{\{0,1\}})^{\otimes T}$

we have, by independence:

$$\begin{aligned} &\text{KL}(\mathbb{P}_i^T \otimes U^T, \mathcal{Q}^T \otimes U^T) \\ &= T \text{KL}(\mathbb{P}_i^{(L_{1T}, \dots, L_{NT})}, \mathcal{Q}^{(L_{1T}, \dots, L_{NT})}) \\ &\leq T \frac{s}{N} \text{KL}\left(\frac{1}{2} - \varepsilon N/s, \frac{1}{2}\right) \leq \frac{4N}{s} \varepsilon^2 T \end{aligned}$$

we will show in a minute

Putting all elements together:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\mathbb{E} \left[\frac{N_i(\tau)}{T} \right] \right] \leq \frac{1}{N} + \sqrt{\frac{4NT\epsilon^2}{s \ln N}}$$

thus $\sup \mathbb{E}[R_T] \geq T\epsilon \left(1 - \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\frac{N_i(\tau)}{T} \right] \right)$

$$\geq T\epsilon \left(\underbrace{1 - \frac{1}{N}}_{\geq \frac{1}{2}} - \sqrt{\frac{4NT\epsilon^2}{s \ln N}} \right) \geq T\epsilon \left(\frac{1}{2} - c\epsilon \right)$$

Optimizing the bound: (maximizing) $\epsilon = 1/4c$ where $c = \sqrt{\frac{4NT}{s \ln N}}$ leads to $T\epsilon/4$, which is the stated bound

But we have to check $\epsilon < s/2N$ (cf first line of the proof):

OK if $T > (N \ln N) / 16s$.

So, it only remains to prove:

$$\left\{ \begin{array}{l} \text{KL}(\mathbb{P}_i^{(L_{1T} \dots L_{NT})}, \mathbb{Q}^{(L_{1T} \dots L_{NT})}) \leq \text{KL}\left(\frac{1}{2} - \frac{\epsilon N}{s}, \frac{1}{2}\right) \\ \text{KL}\left(\frac{1}{2} - \frac{\epsilon N}{s}, \frac{1}{2}\right) \leq \frac{4N^2\epsilon^2}{s^2} \end{array} \right.$$

for $u \in (0, 1/2)$, $\text{KL}\left(\frac{1}{2} - u, \frac{1}{2}\right) = \left(\frac{1}{2} - u\right) \ln(1 - 2u) + \left(\frac{1}{2} + u\right) \ln(1 + 2u) \leq \frac{4u^2}{2} = 2u^2$

$\mathbb{P}_i^{(L_{1T} \dots L_{NT})}$ and $\mathbb{Q}^{(L_{1T} \dots L_{NT})}$ are convex combinations of probability distributions of the following form, indexed by the \emptyset subsets of $\{1, \dots, N\}$ with s elements and up to permutations of the Bernoulli distributions in the products below (which does not change the value of the KL divergence):

[when i is picked]

$\binom{N-1}{s-1}$ distributions of the form $\text{Ber}\left(\frac{1}{2} - \frac{\epsilon N}{s}\right) \otimes \text{Ber}\left(\frac{1}{2}\right)^{\otimes (s-1)} \otimes \mathbb{S}_0^{\otimes (N-s)}$
and $\text{Ber}\left(\frac{1}{2}\right)^{\otimes s} \otimes \mathbb{S}_0^{\otimes (N-s)}$

[when i is not picked]

$\binom{N-1}{s}$ distributions of the form $\text{Ber}\left(\frac{1}{2}\right)^{\otimes s} \otimes \mathbb{S}_0^{\otimes (N-s)}$
and $\text{Ber}\left(\frac{1}{2}\right)^{\otimes s} \otimes \mathbb{S}_0^{\otimes (N-s)}$

Only the first set of distributions contributes to the Kullback-Leibler divergence:

$$\text{KL} \left(\prod_{i=1}^N (L_{NT} \dots L_{NT}), Q (L_{NT} \dots L_{NT}) \right) \leq \underbrace{\frac{\binom{N-1}{s-1}}{\binom{N}{s}}}_{= s/N} \text{KL} \left(\text{Ber} \left(\frac{1}{2} - \frac{\epsilon N}{s} \right) \otimes \text{Ber} \left(\frac{1}{2} \right)^{\otimes (N-1)}, \text{Ber} \left(\frac{1}{2} \right) \otimes \text{Ber} \left(\frac{1}{2} \right)^{\otimes (N-1)} \right)$$

joint by convexity of the KL \nearrow

$$= \text{KL} \left(\frac{1}{2} - \frac{\epsilon N}{s}, \frac{1}{2} \right)$$