

Let's first complete the proof of the Lemma:

[ "Hoeffding-Azuma inequality with a random number of summands" ]

Setting: Probability distributions  $\nu_1, \dots, \nu_K$  over  $[0, 1]$   
with respective expectations  $\mu_1, \dots, \mu_K$

At each round,  $I_t \in \{1, \dots, K\}$  is picked in a  $\sigma(y_1, \dots, y_{t-1})$ -measurable way

then  $y_t$  is drawn independently at random according to  $\nu_{I_t}$ , given  $I_t$

$$\text{i.e.: } y_t | I_t \sim \nu_{I_t}$$

$$\text{We denote } N_a(t) = \sum_{s=1}^t \mathbb{1}_{\{I_s=a\}}$$

and assume that each arm  $a$  was pulled once in the first  $K$  rounds,

so that:

$$N_a(t) \geq 1 \quad \forall t \geq K$$

$$\text{Then, for } t \geq K: \quad \hat{\mu}_{a,t} = \frac{1}{N_a(t)} \sum_{s=1}^t y_s \mathbb{1}_{\{I_s=a\}}$$

$$\text{Lemma: } \forall \delta \in (0, 1), \quad \mathbb{P} \left\{ \mu_a > \hat{\mu}_{a,t} - \sqrt{\frac{\ln(V\delta)}{2N_a(t)}} \right\} \geq 1 - t\delta$$

(and a similar symmetric statement with  $\mu_a < \hat{\mu}_{a,t} + \sqrt{\dots}$ )

The proof will be based on the fact that  $(Z_t)_{t \geq 0}$ , where

$$Z_t = \sum_{s=1}^t (y_s - \mu_a) \mathbb{1}_{\{I_s=a\}}$$

is a martingale w.r.t.  $(\mathcal{F}_t)_{t \geq 0} = (\sigma(y_1, \dots, y_t))_{t \geq 0}$ , which we already proved last

time:

$$\mathbb{E} \left[ (y_t - \mu_a) \mathbb{1}_{\{I_t=a\}} \mid \mathcal{F}_{t-1} \right] = \mathbb{E} \left[ (y_t - \mu_a) \mathbb{1}_{\{I_t=a\}} \mid I_t, y_1, \dots, y_{t-1} \right]$$

where we used the bounded model

$$= \left( \mathbb{E}[y_t \mid I_t, y_1, \dots, y_{t-1}] - \mu_a \right) \mathbb{1}_{\{I_t=a\}} \\ = (\mu_{I_t} - \mu_a) \mathbb{1}_{\{I_t=a\}} = 0 \text{ a.s.}$$

Remark:

How does this bound compare to what the classical version of the Hoeffding-Azuma says?

Martingale increment  $(y_s - \mu_a) \mathbb{1}_{\{I_s=a\}}$  bounded between

$$a_t = -\mu_a \quad \text{and} \quad b_t = 1 - \mu_a$$

so that (actually in the version I stated, I can have  $\leq$  or  $<$ )

$$(b_t - a_t)^2 = 1$$

$$1 - \epsilon \leq \mathbb{P} \left\{ Z_t < \sqrt{\frac{t}{2}} \ln \frac{1}{\epsilon} \right\} = \mathbb{P} \left\{ N_t(t) (\hat{\mu}_t - \mu_a) < \sqrt{\frac{t}{2}} \ln \frac{1}{\epsilon} \right\}$$

$$= \mathbb{P} \left\{ \hat{\mu}_t - \sqrt{\frac{t}{N_t(t)}} \sqrt{\frac{\ln(1/\epsilon)}{2}} < \mu_a \right\}$$

versus the bound of our lemma:  $1 - \epsilon \leq \mathbb{P} \left\{ \hat{\mu}_t - \sqrt{\frac{\ln(1/\epsilon)}{2N_t(t)}} < \mu_a \right\}$

The proposed deviation essentially differ from a  $\sqrt{t/N_t(t)}$  factor, and it is so nice to get rid of it!

Proof: (1) We prove that  $\forall x \in \mathbb{R}, \mathbb{E} \left[ e^{xZ_t - \frac{x^2}{8} N_t(t)} \right] \leq 1$

We do so by showing that  $M_t = \exp \left( xZ_t - \frac{x^2}{8} N_t(t) \right)$  is a supermartingale, so that  $\mathbb{E}[M_t] \leq \mathbb{E}[M_0] = 1$ .

Indeed, by the conditional version of Hoeffding's lemma,

$$\mathbb{E} \left[ e^{x(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \leq e^{x^2/8} \quad \text{a.s.} \quad \left. \begin{array}{l} \text{but we} \\ \text{can do} \\ \text{better!} \end{array} \right\}$$

Since  $\mathcal{I}_t$  and thus also  $\mathbb{1}_{\mathcal{I}_t = a_j}$  are  $\mathcal{F}_{t-1}$ -measurable, we get:

$$\begin{aligned} \mathbb{E} \left[ e^{x(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] &= \mathbb{E} \left[ e^{x(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} (\mathbb{1}_{\mathcal{I}_t = a_j} + \mathbb{1}_{\mathcal{I}_t \neq a_j}) \mid \mathcal{F}_{t-1} \right] \\ &= \mathbb{E} \left[ e^{x(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \mathbb{1}_{\mathcal{I}_t = a_j} + e^0 \mathbb{1}_{\mathcal{I}_t \neq a_j} \\ &\stackrel{\text{given what we had before}}{\leq} e^{x^2/8} \mathbb{1}_{\mathcal{I}_t = a_j} + \mathbb{1}_{\mathcal{I}_t \neq a_j} = \exp \left( \frac{x^2}{8} \mathbb{1}_{\mathcal{I}_t = a_j} \right) \end{aligned}$$

Put differently, 
$$\mathbb{E} \left[ e^{x(Y_t - \mu_a) \mathbb{1}_{\mathcal{I}_t = a_j} - \frac{x^2}{8} \mathbb{1}_{\mathcal{I}_t = a_j}} \mid \mathcal{F}_{t-1} \right] \leq 1$$

which entails that 
$$\exp \left( x \sum_{s=1}^t (Y_s - \mu_a) \mathbb{1}_{\mathcal{I}_s = a_j} - \frac{x^2}{8} \sum_{s=1}^t \mathbb{1}_{\mathcal{I}_s = a_j} \right)$$

$$= \exp \left( xZ_t - \frac{x^2}{8} N_t(t) \right) = M_t$$

is a supermartingale wrt  $\mathcal{F}_t = \sigma(Y_1, \dots, Y_t)$ .

$$(2) \text{ We prove that } \forall \varepsilon > 0, \forall \ell \geq 1, \mathbb{P}\{Z_\ell \geq \varepsilon \text{ and } N_\ell(t) = \ell\} \leq \exp(-2\varepsilon^2/\ell)$$

Indeed, by a Markov-Chernoff bounding,

$$\begin{aligned} \forall x > 0, \quad \mathbb{P}\{Z_\ell \geq \varepsilon \text{ and } N_\ell(t) = \ell\} &\leq e^{-x\varepsilon} \mathbb{E}\left[e^{xZ_\ell} \mathbb{1}_{\{N_\ell(t) = \ell\}}\right] \\ &= e^{-x\varepsilon + \frac{x^2\ell}{3}} \mathbb{E}\left[e^{xZ_\ell - \frac{x^2}{3}N_\ell(t)} \mathbb{1}_{\{N_\ell(t) = \ell\}}\right] \\ &\leq e^{-x\varepsilon + \frac{x^2\ell}{3}} \underbrace{\mathbb{E}\left[e^{xZ_\ell - \frac{x^2}{3}N_\ell(t)}\right]}_{\leq 1 \text{ by (1)}} \end{aligned}$$

Optimizing over  $x > 0$

(take  $x = 4\varepsilon/\ell$ ) yields the claimed bound.

$$(3) \text{ Conclusion: we prove that } \mathbb{P}\left\{\mu_n \leq \hat{\mu}_{\text{opt}} - \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}}\right\} \leq t\delta$$

Indeed, by distinguishing according to the values taken by  $N_n(t)$ :

$$\begin{aligned} &\mathbb{P}\left\{\mu_n \leq \hat{\mu}_{\text{opt}} - \sqrt{\frac{\ln(1/\delta)}{2N_n(t)}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } \mu_n \leq \hat{\mu}_{\text{opt}} - \sqrt{\frac{\ln(1/\delta)}{2\ell}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } \frac{Z_\ell}{N_n(t)} \geq \sqrt{\frac{\ln(1/\delta)}{2\ell}}\right\} \\ &= \sum_{\ell=1}^t \mathbb{P}\left\{N_n(t) = \ell \text{ and } Z_\ell \geq \sqrt{\ell \ln(1/\delta)/2}\right\} \\ &\stackrel{\text{by (2)}}{\leq} \sum_{\ell=1}^t \exp(-2(\ell \ln(1/\delta)/2)/\ell) = t\delta. \end{aligned}$$

(  $\sum_{\ell=1}^{t-K+1}$  would be enough )

Stochastic bandits :What about arms indexed by a continuum?

Setting 1 : Arms indexed by  $x \in A$ , where  $A$  is some possibly large set  
 With each arm  $x \in A$  is associated a probability distribution  $\nu_x$  over  $\mathbb{R}$  s.t.  $E(\nu_x)$  exists  
 At each round, the decision-maker picks  $I_t \in A$ ,  
 gets a reward  $Y_t$  drawn at random according to  $\nu_{I_t}$   
 (given  $I_t$ ); and this is the only feedback she gets.

Definition :  $f: x \in A \mapsto E(\nu_x)$  is the mean-payoff function  
 Regret : 
$$\bar{R}_T = T \sup_{x \in A} f(x) - E\left[\sum_{t=1}^T Y_t\right]$$

Setting 2 : [special case]  $\rightarrow$  Noisy optimization of a function.  
 We fix  $f: A \rightarrow \mathbb{R}$   
 The noise is given by a sequence of iid random variables  $\varepsilon_1, \varepsilon_2, \dots$   
 When  $I_t \in A$  is picked,  $Y_t = f(I_t) + \varepsilon_t$   
 $\hookrightarrow$  Special case of setting #1 where  $\nu_x$  is the distribution of  $f(x) + \varepsilon_1$  (all these distributions have the same shape given by the common distribution of the  $\varepsilon_j$ )

We of course need conditions for the regret to be minimized.

Definition : Let  $\mathcal{F}$  be a set of possible bandit problems  $\mathcal{F} = (\nu_x)_{x \in A}$   
 The regret can be controlled (in a non-uniform way) against  $\mathcal{F}$  if :

we also say that  $(A, \mathcal{F})$  is tractable

there exists a strategy s.t.  $\forall \mathcal{F} \in \mathcal{F}, \bar{R}_T = o(T)$ .



Ex:  $A = \{1, \dots, K\}$  and  $\mathcal{F} = (\mathcal{P}([0,1]))^K$ , the set of all  $K$ -tuples of probability distributions over  $[0,1]$   
 the case of finitely many arms with bounded distributions  
 → UCB does the job.

Counter-example:  $A = [0,1]$  and  $\mathcal{F} = (\mathcal{P}([0,1]))^{[0,1]}$   
 illustrating that continuity is a minimal requirement.  
 all bandit problems  $(\nu_x)_{x \in [0,1]}$  with distributions  $\nu_x$  having support  $[0,1]$ .

Includ: Consider  $(\delta_0)_{x \in [0,1]}$  the bandit problem in which each arm  $x$  is associated with the Dirac mass on 0

Fix any strategy: it gets  $Y_t = 0$   $\forall t$  and uses a sequence of (possibly) random choices  $I_t, t \geq 1$   
 Since probability distributions can only have at most countably many atoms,  
 $\mathcal{Y} = \{x \in [0,1] : \exists t \mid \mathbb{P}[I_t = x] > 0 \text{ under } (\delta_0)_{x \in [0,1]}\}$  is countable. In particular,  $[0,1] \setminus \mathcal{Y}$  is non empty.

But the strategy behaves the same under the problem  $(\nu'_x)_{x \in [0,1]}$  in which  $\begin{cases} \nu'_x = \delta_0 & \forall x \neq x_0 \\ \nu'_{x_0} = \delta_1 & \text{for one fixed } x_0 \in [0,1] \setminus \mathcal{Y} \end{cases}$

With probability 1, it thus never hits  $x_0$ .

Therefore,  $Y_t = 0$  a.s.  $\forall t$  and  $\bar{R}_T = T - \mathbb{E}[\sum_{t=1}^T Y_t] = T$

Actually, continuity is sufficient for the regret to be controlled, as long as  $A$  is not too large.

Theorem: Let  $A$  be a compact metric space and let  $\mathcal{F}$  be the set of bandit problems  $(\nu_x)_{x \in A}$

- with: →  $\forall x, \nu_x$  is a distribution over  $[0,1]$
- a continuous mean-payoff function  $f, x \mapsto \mathbb{E}[\nu_x]$

The regret can be controlled against  $\mathcal{F}^{\text{cont}}$  if and only if  $A$  is separable.

Corollary. Let  $\mathcal{F}^{\text{all}}$  be the family of all bandit models  $(\vec{\mu}_x)_{x \in A}$  with distributions  $\nu_x$  over  $[0,1]$ . Then the regret against  $\mathcal{F}^{\text{all}}$  can be controlled if and only if  $A$  is at most countable.

Before we prove these facts, consider the following more concrete example, in which, by strengthening the regularity requirement on the mean-payoff function, we can even get rates.

Exercise: (Lipschitz bandits) Let  $A = [0,1]$  and let  $\mathcal{F}^{\text{lip}}$  be the family of bandit models  $(\vec{\mu}_x)_{x \in [0,1]}$  with distributions  $\nu_x$  over  $[0,1]$  and with mean-payoff functions that are Lipschitz.

Exhibit a strategy based on UCB + a sequence of discretizations of  $[0,1]$  into  $K$  bins (to be refined over time) such that:

Hint:

First, prove a performance bound by splitting  $[0,1]$  into  $[(i-1)/K, i/K]$  with  $i=1, \dots, K$  for a fixed  $K$ , where each bin  $[(i-1)/K, i/K]$  plays the role of an  $i$  in a bandit problem with finitely many arms. Then discuss how to pick  $K$  over the time, as we do in the next proof.

$\forall \epsilon \in \mathcal{F}^{\text{lip}}$ ,

$$\bar{R}_T \leq (3L + 6\sqrt{8 \ln T + 2})T^{2/3} + 2$$

where  $L$  is the Lipschitz constant of the mean-payoff function of  $\epsilon$ .

Proof of the Corollary:

We endow  $A$  with the discrete topology, i.e., choose the distance

$$d(x, y) = \mathbb{1}_{\{x \neq y\}}. \quad \text{Then:}$$

1. All applications  $f: A \rightarrow \mathbb{R}$  are continuous
2.  $A$  is separable if and only if  $A$  is at most countable.

Proof of the Theorem:

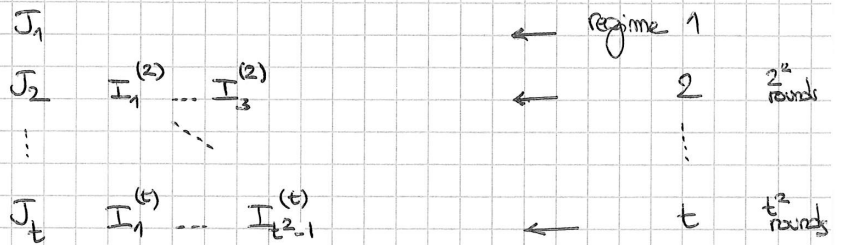
It relies on the possibility or impossibility of uniform exploration of the arms.

1) If  $A$  is separable: let  $(x_n)_{n \in \mathbb{N}}$  be a collection of points in  $A$  that is dense

In particular, the probability distribution  $\mu = \sum_{n \geq 0} \frac{1}{2^{n+1}} \delta_{x_n}$  is such that  $\mu(V) > 0$  for all open sets  $V \subset A$ .

We pick elements  $J_1, J_2, I_1^{(2)}, \dots, J_t, I_1^{(t)}, \dots, I_{t-1}^{(t)}$  as follows:

Actually we could take  $J_m = x_n$  if it would work, this extra randomization is not needed...



where  $\{$  the  $J_s$  are drawn at random according to  $\mu$   
 the  $I_s^{(r)}$ ,  $1 \leq s \leq r^2 - 1$ , follow from the UCS strategy with arms  $J_1, \dots, J_r$

In regime  $r$ :

$$r^2 \max_{s \leq r} \mu_{J_s} - \mathbb{E} \left[ \sum_{s=1}^{r^2} Y_{S_{r^2+s}} \right]$$

regime  $r$  starts at time  $S_{r-1} + 1$  where  $S_0 = 1, S_1 = 1+2, \dots, S_{r-1} = \frac{r(r-1)}{2}$

we have it first for  $\mathbb{E}[Y_{J_r}]$  and then of course for  $\mathbb{E}[Y]$  by tower rule

$$\leq 1 + c \sqrt{r^3 \ln r}$$

↑  
for  $J_r$

↑  
well-chosen numerical constant

↑  
distribution-free regret bound for UCS on  $r^2 - 1$  steps with  $r$  arms (we saw this bound as an exercise)

Let  $\epsilon > 0$ , let  $\tilde{r}_\epsilon$  the first (random) time when  $\mu_{J_r} = f(J_r) \geq \sup_{z \in A} f(z) - \epsilon$

We have  $\tilde{r}_\epsilon < +\infty$  a.s. because:

- by continuity of  $f$ , there exists an open set  $V_\epsilon$  with  $\forall x \in V_\epsilon, f(x) \geq \sup_A f - \epsilon$ ;

- we have  $\tilde{r}_\varepsilon \leq \inf \{ r \geq 1 : J_r \in V_\varepsilon \} < +\infty$  a.s.  
 as this random variable follows a geometric distribution with parameter  $\nu(V_\varepsilon) > 0$ .

For  $r \geq \tilde{r}_\varepsilon$ ,  $\max_{S \leq r} \mu_{J_S} + \varepsilon \geq \sup_A f$

So that 
$$\bar{R}_T = T \sup_A f - \mathbb{E} \left[ \sum_{t=1}^T J_t \right]$$

$$\leq \sum_{r=1}^{\tilde{r}_T-1} r^2 + T\varepsilon + \sum_{r=\tilde{r}_T}^{r_T-1} (1 + c\sqrt{r^3 \ln T}) + r_T^2$$

↑ lengths of regimes  $r \leq \tilde{r}_T - 1$   $< +\infty$  a.s.  
 ↑ regime  $r_T$  may be incomplete  
 where  $r_T$  is such that time  $T$  belongs to regime  $r_T$ :  
 we have  $r_T^3$  of order  $T$   
 i.e.  $r_T$  of order  $T^{1/3}$

and

$$\sum_{r \leq r_T} (1 + c r^{3/2} \sqrt{\ln r}) \leq \sum_{r \leq r_T} (1 + c r^{3/2} \sqrt{\ln r_T}) = O(r_T^{5/2} \sqrt{\ln r_T}) = O(T^{5/6} \sqrt{\ln T})$$

Thus,  $\limsup_{T \rightarrow +\infty} \frac{\bar{R}_T}{T} \leq \varepsilon$  a.s.

but since  $\bar{R}_T$  is a deterministic quantity and this is true  $\forall \varepsilon > 0$ , we have

$$\lim_{T \rightarrow +\infty} \frac{\bar{R}_T}{T} = 0 \text{ as requested.}$$



2) If  $A$  is not separable:

\* We use the following characterization of separability (which relies on Zorn's lemma):

A metric space  $X$  is separable if and only if it contains no uncountable subset  $\mathcal{D}$  s.t.  

$$\rho = \inf \{ d(x, y) : x, y \in \mathcal{D} \} > 0.$$

In particular, if  $A$  is not separable, there exist an uncountable subset  $\mathcal{D} \subset A$  and  $\rho > 0$  such that the balls  $B(a, \rho/2)$ , with  $a \in \mathcal{D}$ , are all disjoint.

$\Rightarrow$  No probability distribution over  $A$  can give a positive mass to all these balls.

\* We consider the bandit models  $\nu^{(a)}$  inducing mean-payoff functions  $f^{(a)} : x \in A \mapsto (1 - \frac{d(x, a)}{\rho/2})^+$ ; in particular,  $\nu^{(a)} = \delta_a$  for  $x \notin B(a, \rho/2)$ .  
 $\uparrow$   $f^{(a)}$  is indeed continuous.

We proceed as in the example showing the necessity of continuity when  $A = [0, 1]$  and consider the bandit model  $(\delta_a)_{a \in A}$ , as well as any strategy and the laws induced by the  $\mathcal{F}_t$  under this model: let  $\nu_t$  be the law of  $\mathcal{F}_t$  under  $(\delta_a)_{a \in A}$  and let  $d = \sum_{t \geq 1} \frac{1}{2^t} \nu_t$ .

There exists  $a \in A$  s.t.  $d(B(a, \rho/2)) = 0$ , that is, s.t.,  $\forall t \geq 1$ ,  $\mathbb{P}(\mathcal{F}_t \in B(a, \rho/2) \text{ under } (\delta_a)_{a \in A}) = 0$ .

The considered strategy is therefore such that the  $\mathcal{F}_t$  have the same distribution under  $(\delta_a)_{a \in A}$  and  $\nu^{(a)}$ . In particular,

$\mathbb{E}[\sum_{t=1}^T Y_t] = 0$  in both cases, but in the latter case,  $\sup f^{(a)} = 1$ , so that  $\bar{R}_T = T$  against  $\nu^{(a)}$ . The regret is not controlled against  $\nu^{(a)} \in \bar{\mathcal{F}}^{\text{cont}}$ .

Overview of the next steps: Fix a model  $\mathcal{D}$ , known to the decision-maker, ie, a collection of probability distributions over  $\mathcal{R}$  with an expectation.

Assume that  $y_1^*, \dots, y_k^*$  are unknown but that the decision-maker knows  $y_j \in \mathcal{D}$  if  $y_j$ .

What are the best bounds on  $\bar{R}_T = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T Y_t \right]$ ?

We will show matching upper and lower bounds (with associated strategies):

$\bar{R}_T$  is at best of the order of  $\left( \sum_{a: \Delta_a > 0} \frac{\Delta_a}{K_{\text{KL}}(\bar{y}_a, \mu^*, \mathcal{D})} \right) \ln T$

where

$$K_{\text{KL}}(\bar{y}_a, \mu^*, \mathcal{D}) = \inf \left\{ KL(\bar{y}_a, \bar{y}_a^*) : \begin{array}{l} \bar{y}_a^* \in \mathcal{D} \\ \mathbb{E}(\bar{y}_a^*) > \mu^* \end{array} \right\}$$

We will do so by

- proving a universal lower bound

Kullback-Leibler divergence

expectation of  $\bar{y}_a$

- exhibiting a strategy, called KL-UCB, to achieve the bound. ← if time permits (I'm not sure we will have time to do so...)

\* Part \* before we do that, I guess that some reminder of basic and non-basic results about KL divergences would be needed!

## The Kullback-Leibler divergence: definition and basic properties.

Definition (intrinsic): Let  $\mathbb{P}, \mathbb{Q}$  be two probability measures over  $(\Omega, \mathcal{F})$

$$KL(\mathbb{P}, \mathbb{Q}) = \begin{cases} +\infty & \text{if } \mathbb{P} \text{ is not absolutely continuous w.r.t } \mathbb{Q} \\ \int_{\Omega} \left( \frac{d\mathbb{P}}{d\mathbb{Q}} \ln \frac{d\mathbb{P}}{d\mathbb{Q}} \right) d\mathbb{Q} = \int_{\Omega} \left( \ln \frac{d\mathbb{P}}{d\mathbb{Q}} \right) d\mathbb{P} & \text{if } \mathbb{P} \ll \mathbb{Q} \end{cases}$$

### Basic facts:

- Existence of the defining integral when  $\mathbb{P} \ll \mathbb{Q}$ : because  $\psi: x \mapsto x \ln x$  is bounded from below on  $[0, +\infty)$

- $KL(\mathbb{P}, \mathbb{Q}) \geq 0$  and  $KL(\mathbb{P}, \mathbb{Q}) = 0$  if and only if  $\mathbb{P} = \mathbb{Q}$ :

It suffices to consider the case  $\mathbb{P} \ll \mathbb{Q}$ : because  $\psi$  is strictly convex, Jensen's inequality indicates that

$$KL(\mathbb{P}, \mathbb{Q}) = \int_{\Omega} \psi\left(\frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{Q} \geq \psi\left(\underbrace{\int_{\Omega} \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{Q}}_{=1}\right) = 0$$

with equality if and only if  $\frac{d\mathbb{P}}{d\mathbb{Q}}$  is  $\mathbb{Q}$ -a.s. constant, i.e.,  $\mathbb{P} = \mathbb{Q}$

Exercise: A useful rewriting. Prove the following result:

Assume  $\mathbb{P} \ll \mathbb{Q}$  and let  $\nu$  be any probability measure over  $(\Omega, \mathcal{F})$

such that  $\mathbb{P} \ll \nu$  and  $\mathbb{Q} \ll \nu$ . Denote  $f = \frac{d\mathbb{P}}{d\nu}$  and  $g = \frac{d\mathbb{Q}}{d\nu}$ .

Then:

$$\begin{aligned} KL(\mathbb{P}, \mathbb{Q}) &= \int_{\Omega} \frac{f}{g} \ln\left(\frac{f}{g}\right) g d\nu \\ &= \int_{\Omega} \ln\left(\frac{f}{g}\right) f d\nu \end{aligned}$$

Beware: with the usual measure-theoretic conventions, if  $x \neq 0$  and  $y = 0$ , then  $x \neq y \frac{x}{y}$   $\hookrightarrow$  you therefore need to proceed with care!

Lemma (Contraction of entropy; also known as data-processing inequality):

Let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$

Let  $X: (\Omega, \mathcal{F}) \rightarrow (\Omega', \mathcal{F}')$  be any random variable

Denote by  $P^X$  and  $Q^X$  the laws of  $X$  under  $P$  and  $Q$ .

Then:

$$KL(P^X, Q^X) \leq KL(P, Q).$$

Proof: We may assume that  $P \ll Q$ , otherwise  $KL(P, Q) = +\infty$  and the inequality is true. We show that we then have

$$P^X \ll Q^X, \quad \text{with} \quad \frac{dP^X}{dQ^X} = \mathbb{E}_Q \left[ \frac{dP}{dQ} \mid X = \cdot \right] \stackrel{\text{not.}}{=} \gamma$$

$$\text{ie } \gamma(x) = \mathbb{E}_Q \left[ \frac{dP}{dQ} \mid X \right].$$

Indeed, for all  $B \in \mathcal{F}'$ :

$$P^X(B) = P\{X \in B\} = \int_{\Omega} \mathbb{1}_B(X) \frac{dP}{dQ} dQ \stackrel{\text{tower rule}}{=} \int_{\Omega} \mathbb{1}_B(X) \mathbb{E}_Q \left[ \frac{dP}{dQ} \mid X \right] dQ$$

$$\stackrel{\text{not.}}{=} \int_{\Omega} \mathbb{1}_B(X) \gamma(x) dQ = \int_{\Omega'} \mathbb{1}_B \gamma dQ^X.$$

Therefore,

$$KL(P^X, Q^X) = \int_{\Omega'} \gamma \ln \gamma dQ^X = \int_{\Omega} \gamma(x) \ln \gamma(x) dQ$$

$$= \int_{\Omega} \left( \mathbb{E}_Q \left[ \frac{dP}{dQ} \mid X \right] \ln \mathbb{E}_Q \left[ \frac{dP}{dQ} \mid X \right] \right) dQ \quad \left. \begin{array}{l} \text{definition} \\ \text{of } \gamma \end{array} \right\}$$

$$\leq \int_{\Omega} \mathbb{E}_Q \left[ \frac{dP}{dQ} \ln \frac{dP}{dQ} \mid X \right] dQ \quad \left. \begin{array}{l} \text{conditional} \\ \text{version of} \\ \text{Jensen's inequality} \end{array} \right\}$$

$$\stackrel{\text{tower rule}}{\rightarrow} = \int_{\Omega} \left( \frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ = KL(P, Q)$$



References: • The proof above is due to Ali and Silvey (1966), but it's far from being well-known!

- Typical proofs in the more recent literature:
  - either focus on the discrete case (Cover and Thomas, 2006)
  - or use the duality / variational formula for the KL (Massart, 2007; Boucheron, Lugosi, Massart, 2013)

• The joint convexity of KL, which we discuss below, is typically proved in a tedious way, relying on the rewriting of Exercise 1 and the joint convexity of  $(x, y) \in [0, +\infty)^2 \mapsto \left(\frac{x}{y} \ln \frac{x}{y}\right) y$

We may see it instead as a consequence of the data-processing inequality:

Corollary (joint convexity of KL): For all probability distributions  $\mathbb{P}_1, \mathbb{P}_2$  and  $\mathbb{Q}_1, \mathbb{Q}_2$  over the same measurable space  $(\Omega, \mathcal{F})$ , and all  $d \in (0, 1)$ ,

$$KL((1-d)\mathbb{P}_1 + d\mathbb{P}_2, (1-d)\mathbb{Q}_1 + d\mathbb{Q}_2) \leq (1-d) KL(\mathbb{P}_1, \mathbb{Q}_1) + d KL(\mathbb{P}_2, \mathbb{Q}_2)$$

Proof: We augment  $(\Omega, \mathcal{F})$  into  $(\Omega \times \{1, 2\}, \mathcal{F}')$  where  $\mathcal{F}' = \mathcal{F} \otimes \{\emptyset, \{1\}, \{2\}, \{1, 2\}\}$

We define the random pair  $(X, J)$  by the projections

$$X: (\omega, j) \mapsto \omega \quad \text{and} \quad J: (\omega, j) \mapsto j$$

Let  $\mathbb{P}$  be a probability measure on  $(\Omega \times \{1, 2\}, \mathcal{F}')$  such that

$$\begin{cases} J \sim 1 + \text{Ber}(d) \\ X | J=j \sim \mathbb{P}_j \end{cases} \quad (\text{and a similar definition for } \mathbb{Q} \text{ based on } \mathbb{Q}_1, \mathbb{Q}_2)$$

that is,  $\forall j \in \{1, 2\} \quad \forall A \in \mathcal{F} \quad \mathbb{P}(A \times \{j\}) = \left( (1-d) \mathbb{1}_{\mathbb{P}_j} + d \mathbb{1}_{\mathbb{P}_{j=2}} \right) \mathbb{P}_j(A)$

$$\begin{aligned} P^X &= (1-d)P_1 + dP_2 \\ Q^X &= (1-d)Q_1 + dQ_2 \end{aligned}$$

and (as we prove below)  $KL(P^X, Q^X) = (1-d)KL(P_1, Q_1) + dKL(P_2, Q_2)$   
 so that the result follows from the data-processing inequality.

Indeed: we may assume with no loss of generality, given  $d \in (0,1)$ , that  $P_1 \ll Q_1$  and  $P_2 \ll Q_2$ , so that  $P \ll Q$  with

$$\frac{dP}{dQ}(w, j) = \mathbb{1}_{\{j=1\}} \frac{dP_1}{dQ_1}(w) + \mathbb{1}_{\{j=2\}} \frac{dP_2}{dQ_2}(w)$$

This entails that

$$\begin{aligned} KL(P, Q) &= \int_{\Omega \times \{1,2\}} \left( \frac{dP}{dQ}(w, j) \ln \frac{dP}{dQ}(w, j) \right) dQ(w, j) \\ &= \int_{\Omega \times \{1\}} \left( \frac{dP}{dQ}(w, 1) \ln \frac{dP}{dQ}(w, 1) \right) \mathbb{1}_{\Omega \times \{1\}}(w, j) dQ(w, j) \\ &\quad + \int_{\Omega \times \{2\}} \left( \frac{dP}{dQ}(w, 2) \ln \frac{dP}{dQ}(w, 2) \right) \mathbb{1}_{\Omega \times \{2\}}(w, j) dQ(w, j) \\ &= \int_{\Omega} \left( \frac{dP_1}{dQ_1}(w) \ln \frac{dP_1}{dQ_1}(w) \right) dQ_1(w) + \int \dots \\ &= KL(P_1, Q_1) + KL(P_2, Q_2) \end{aligned}$$

*we just use that for  $f \geq a$  constant,  $\int f d\mu = \int f \mathbb{1}_A d\mu + \int f \mathbb{1}_{A^c} d\mu$  whether  $f$  is integrable or not*

*on  $\Omega \times \{1\}$ ,  $dQ$  is  $d dQ_1$*

KL for product measures. ( $\leftrightarrow$  The independent case)

Proposition: Let  $(\Omega, \mathcal{F})$  and  $(\Omega', \mathcal{F}')$  be two measurable spaces,  
 let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$   
 $P', Q'$  over  $(\Omega', \mathcal{F}')$

and denote by  $P \otimes P'$  and  $Q \otimes Q'$  the product distributions over  $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$ . Then:

$$KL(P \otimes P', Q \otimes Q') = KL(P, Q) + KL(P', Q')$$

Proof: We have  $P \otimes P' \ll Q \otimes Q' \iff [P \ll Q \text{ and } P' \ll Q']$

so we can assume that all  $\ll$  statements hold, and then

$$\frac{d(P \otimes P')}{d(Q \otimes Q')} = \frac{dP}{dQ} \frac{dP'}{dQ'}$$

(this is a fundamental result in measure theory and one of the best characterizations of independence!).

Therefore,

$$KL(P \otimes P', Q \otimes Q') = \int_{\Omega \times \Omega'} \left( \frac{dP}{dQ} \frac{dP'}{dQ'} \ln \left( \frac{dP}{dQ} \frac{dP'}{dQ'} \right) \right) d(Q \otimes Q')$$

We use that if  $f, g$  are  $\geq$  a constant, then  $\int (f+g) d\mu = \int f d\mu + \int g d\mu$

$$= \int_{\Omega} \left( \int_{\Omega'} \left( \frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ' \right) \frac{dP'}{dQ'} dQ + \text{similar term with } \ln \frac{dP'}{dQ'}$$

$\underbrace{\hspace{10em}}_{= KL(P, Q)} \quad \quad \quad \underbrace{\hspace{10em}}_{= KL(P', Q')}$

$\underbrace{\hspace{15em}}_{= KL(P, Q)}$

here we apply Tonelli's theorem (again because  $x \mapsto x \ln x$  is lower bounded)

Consequence (Garivier, Néron, Stoltz, 2016):

Data-processing inequality with expectations of random variables

Corollary: Let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$

Let  $X: (\Omega, \mathcal{F}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$  be any  $[0, 1]$ -valued random variable

Then, denoting by  $E_P[X]$  and  $E_Q[X]$  the respective expectations of  $X$  under  $P$  and  $Q$ , we have:

$$E_P[X] \ln \frac{E_P[X]}{E_Q[X]} + (1 - E_P[X]) \ln \frac{1 - E_P[X]}{1 - E_Q[X]} = KL(\text{Ber}(E_P[X]), \text{Ber}(E_Q[X])) \leq KL(P, Q)$$

Proof: We denote by  $m$  the Lebesgue measure over  $[0, 1]$  and augment the underlying measurable space into  $(\Omega \times [0, 1], \mathcal{F} \otimes \mathcal{B}([0, 1]))$ , over which we consider the product-distributions  $P \otimes m$  and  $Q \otimes m$ .

For any event  $E \in \mathcal{F} \otimes \mathcal{B}([0, 1])$ , we have, by the data-processing inequality:

$$\begin{aligned}
 \text{KL}\left(\underbrace{(\mathbb{P} \otimes \eta)}^{\text{Ber}(\mathbb{P} \otimes \eta(E))}, \underbrace{(\mathbb{Q} \otimes \eta)}^{\text{Ber}(\mathbb{Q} \otimes \eta(E))}\right) &\leq \text{KL}(\mathbb{P} \otimes \eta, \mathbb{Q} \otimes \eta) \\
 &= \text{KL}(\mathbb{P}, \mathbb{Q}) + \text{KL}(\eta, \eta) \\
 &\stackrel{\substack{\uparrow \\ \text{of product} \\ \text{distributions}}}{=} \text{KL}(\mathbb{P}, \mathbb{Q})
 \end{aligned}$$

Thus:  $\text{KL}(\text{Ber}(\mathbb{P} \otimes \eta(E)), \text{Ber}(\mathbb{Q} \otimes \eta(E))) \leq \text{KL}(\mathbb{P}, \mathbb{Q})$

The proof is concluded by picking  $E \in \mathcal{F} \otimes \mathcal{B}([a,1])$  such that  $\mathbb{P} \otimes \eta(E) = \mathbb{E}_{\mathbb{P}}[x]$  and  $\mathbb{Q} \otimes \eta(E) = \mathbb{E}_{\mathbb{Q}}[x]$

Namely,  $E = \{(\omega, x) \in \Omega \times [a,1] : x \leq X(\omega)\}$

By Tonelli's theorem:

$$\begin{aligned}
 \mathbb{P} \otimes \eta(E) &= \int_{\Omega} \left( \int_{[a,1]} \mathbb{1}_{\{x \leq X(\omega)\}} d\eta(x) \right) d\mathbb{P}(\omega) \\
 &= \int_{\Omega} X(\omega) d\mathbb{P}(\omega) = \mathbb{E}_{\mathbb{P}}[x]
 \end{aligned}$$

and a similar equality for  $\mathbb{Q} \otimes \eta(E)$ .