

**Part 1: The Hoeffding–Azuma inequality
with a random number of summand**

Let's first complete the proof of the Lemma:

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

Setting: Probability distributions $\tilde{v}_1, \dots, \tilde{v}_K$ over $[0,1]$
with respective expectations μ_1, \dots, μ_K

At each round, $I_t \in \{1, \dots, K\}$ is picked in a $\sigma(\tilde{y}_1, \dots, \tilde{y}_{t-1})$ -measurable way

then y_t is drawn independently at random according to \tilde{v}_{I_t} , given I_t
i.e.: $y_t | I_t \sim \tilde{v}_{I_t}$

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

Then, for $t \geq K$:

$$\hat{\mu}_a = \frac{1}{N_a(t)} \sum_{s=1}^t y_s \mathbb{1}_{\{I_s=a\}}$$

Lemma: $\forall \delta \in (0,1)$, $P\left\{ \mu_a > \hat{\mu}_a + \sqrt{\frac{\ln(1/\delta)}{2N_a(t)}} \right\} \geq 1 - \delta$
(and a similar symmetric statement with $\mu_a < \hat{\mu}_a + \sqrt{\cdot}$)

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. $(\mathcal{F}_t) = (\sigma(y_1, \dots, y_t))_{t \geq 0}$, which we already proved last time:

$$\begin{aligned} \mathbb{E}\left[(y_t - \mu_a) \mathbb{1}_{\{I_t=a\}} \mid \mathcal{F}_{t-1} \right] &= \mathbb{E}\left[(y_t - \mu_a) \mathbb{1}_{\{I_t=a\}} \mid I_t, \mathcal{F}_{t-1} \right] \\ &= (\mathbb{E}[y_t \mid I_t, \mathcal{F}_{t-1}] - \mu_a) \mathbb{1}_{\{I_t=a\}}. \end{aligned}$$

where we used the bandit model

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

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

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

so that

(actually in the written I stated, I can have \leq or $<$)

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

$$1-tS \leq \mathbb{P}\left\{ Z_t < \sqrt{\frac{t}{2} \ln \frac{1}{tS}} \right\} = \mathbb{P}\left\{ N_\alpha(t) (\hat{\mu}_{at} - \mu_\alpha) < \sqrt{\frac{t}{2} \ln \frac{1}{tS}} \right\}$$

$$= \mathbb{P}\left\{ \hat{\mu}_{at} - \sqrt{\frac{t}{N_\alpha(t)}} \sqrt{\frac{\ln(1/tS)}{2N_\alpha(t)}} < \mu_\alpha \right\}$$

versus the bound of our lemma: $1-tS \leq \mathbb{P}\left\{ \hat{\mu}_{at} - \sqrt{\frac{\ln(1/tS)}{2N_\alpha(t)}} < \mu_\alpha \right\}$

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

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

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

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

$$\mathbb{E}[e^{x(Y_t - \mu_\alpha) \mathbf{1}_{\{I_t=a\}}} | \mathcal{F}_{t-1}] \leq e^{x^2/8} \text{ a.s. } \quad \begin{matrix} \text{but we} \\ \text{can do better!} \end{matrix}$$

Since I_t and thus also $\mathbf{1}_{\{I_t=a\}}$ are \mathcal{F}_{t-1} -measurable, we get:

$$\begin{aligned} \mathbb{E}[e^{x(Y_t - \mu_\alpha) \mathbf{1}_{\{I_t=a\}}} | \mathcal{F}_{t-1}] &= \mathbb{E}[e^{x(Y_t - \mu_\alpha) \mathbf{1}_{\{I_t=a\}} (\mathbf{1}_{\{I_t=a\}} + \mathbf{1}_{\{I_t \neq a\}})} | \mathcal{F}_{t-1}] \\ &= \mathbb{E}[e^{x(Y_t - \mu_\alpha) \mathbf{1}_{\{I_t=a\}}} | \mathcal{F}_{t-1}] \mathbf{1}_{\{I_t=a\}} + e^0 \mathbf{1}_{\{I_t \neq a\}} \\ &\stackrel{\text{given what we had before}}{\leq} e^{x^2/8} \mathbf{1}_{\{I_t=a\}} + \mathbf{1}_{\{I_t \neq a\}} = \exp\left(\frac{x^2}{8} \mathbf{1}_{\{I_t=a\}}\right) \end{aligned}$$

Put differently, $\mathbb{E}[e^{x(Y_t - \mu_\alpha) \mathbf{1}_{\{I_t=a\}} - \frac{x^2}{8} \mathbf{1}_{\{I_t=a\}}} | \mathcal{F}_{t-1}] \leq 1$

which entails that

$$\begin{aligned} &\exp\left(x \sum_{s=1}^t (Y_s - \mu_\alpha) \mathbf{1}_{\{I_s=a\}} - \frac{x^2}{8} \sum_{s=1}^t \mathbf{1}_{\{I_s=a\}}\right) \\ &= \exp\left(xZ_t - \frac{x^2}{8} N_\alpha(t)\right) = M_t \end{aligned}$$

is a supermartingale w.r.t $\mathcal{F}_t = \sigma(Y_1 \dots Y_t)$.

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

Indeed, by a Markov-Chernoff bounding,

$$\begin{aligned} \forall x > 0, \quad \mathbb{P}\{Z_t \geq \varepsilon \text{ and } N_a(t) = l\} &\leq e^{-x\varepsilon} \mathbb{E}[e^{xZ_t} \mathbb{1}_{\{N_a(t) = l\}}] \\ &= e^{-x\varepsilon + \frac{x^2l}{8}} \mathbb{E}[e^{xZ_t - \frac{x^2N_a(t)}{8}} \mathbb{1}_{\{N_a(t) = l\}}] \\ &\leq e^{-x\varepsilon + \frac{x^2l}{8}} \underbrace{\mathbb{E}[e^{xZ_t - \frac{x^2N_a(t)}{8}}]}_{\leq 1 \text{ by (1)}} \end{aligned}$$

Optimizing over $x > 0$

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

(3) Conclusion: we prove that $\mathbb{P}\{\mu_a \leq \hat{\mu}_a - \sqrt{\frac{\ln(1/\delta)}{2N_a(t)}}\} \leq \delta$

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

$$\begin{aligned} &\mathbb{P}\{\mu_a \leq \hat{\mu}_a - \sqrt{\frac{\ln(1/\delta)}{2N_a(t)}}\} \\ &= \sum_{l=1}^t \mathbb{P}\{N_a(t) = l \text{ and } \mu_a \leq \hat{\mu}_a - \sqrt{\frac{\ln(1/\delta)}{2l}}\} \\ &= \sum_{l=1}^t \mathbb{P}\{N_a(t) = l \text{ and } \frac{Z_t}{N_a(t)} \geq \sqrt{\frac{\ln(1/\delta)}{2l}}\} \\ &= \sum_{l=1}^t \mathbb{P}\{N_a(t) = l \text{ and } Z_t \geq \sqrt{l \ln(1/\delta)/2}\} \\ &\stackrel{\text{by (2)}}{\leq} \sum_{l=1}^t \exp(-2(l \ln(1/\delta)/2)/l) = t\delta. \end{aligned}$$

Two notes on this proof:

- * We noted last week that the conditional version of Hoeffding's lemma (based on Proof #2) could be generalized into

X bounded random variable, U, V two \mathcal{G}_j -measurable random variables
with $U \leq X \leq V$

then $\forall s \in \mathbb{R}$,

$$\ln \mathbb{E}[e^{sx} | \mathcal{G}_j] \leq s \mathbb{E}[x | \mathcal{G}_j] + \frac{s^2}{8}(V-U)^2$$

(Because of the proof technique relying on weights with denominators $\sqrt{V-U}$, it is safer to assume $V-U > 0$ as, eg, by replacing U and V by $U-\varepsilon$ and $V+\varepsilon$ if needed and letting $\varepsilon \downarrow 0$; so at the end of the day we may drop the $V-U > 0$ as condition.)

* This extension may be applied to

$$X_t = (\beta_t - \mu_a) \mathbb{1}_{\{I_t=a\}}$$

$$G_j = \mathcal{F}_{t-1}$$

$$U_t = -\mu_a \mathbb{1}_{\{I_t=a\}}$$

$$V_t = (1-\mu_a) \mathbb{1}_{\{I_t=a\}}$$

and directly entails

$$\mathbb{E}[e^{x(\beta_t - \mu_a) \mathbb{1}_{\{I_t=a\}}} | \mathcal{F}_{t-1}] \leq \exp\left(\frac{s^2}{8} \mathbb{1}_{\{I_t=a\}}\right)$$

without the need for the

$$1 = \mathbb{1}_{\{I_t=a\}} + \mathbb{1}_{\{I_t \neq a\}}$$

trick used in Step (1).

* The question is:

Don't we have a generalized version of the Hoeffding-Azuma inequality with such predictable ranges $V_t - U_t$?

Yes we do have something in terms of constant upper bounds

$$V_t - U_t \leq \Delta_t \in \mathbb{R} \text{ as}$$

But $V_t - U_t = \mathbb{1}_{\{I_t=a\}}$ can only be bounded by $\Delta_t = 1$

So I think that Steps (2) and (3) are needed

Part 2: Stochastic bandits with continuously many arms

Stochastic bandits:What about arms indexed by a continuum?Setting 1:Arms indexed by $x \in A$, where A is some possibly large setWith each arm $x \in A$ is associated a probability distribution ν_x over \mathbb{R} s.t. $E(\nu_x)$ existsAt each round, the decision-maker picks $I_t \in A$,gets a reward y_t drawn at random according to ν_{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, $y_t = f(I_t) + \varepsilon_t$

↳ Special case of setting #1 where ν_x is the distribution of $f(x) + \varepsilon_1$ (all these distributions have the same shape given by the common distribution of the ε_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. $\forall \mathcal{F}, \bar{R}_T = o(T)$.

Ex: $A = \{1, \dots, K\}$ and $\mathcal{F} = (\mathbb{P}(\mathcal{Q}_1))^K$, the set of all K -tuples of probability distributions over \mathcal{Q}_1 .
 the case of
 finitely many arms
 with bounded distributions
 → UCB does the job.

Counter-example: $A = \mathcal{Q}_1$ and $\mathcal{F} = (\mathbb{P}(\mathcal{Q}_1))^{\mathcal{Q}_1}$

↑
 illustrating that
 continuity is a minimal
 requirement.
 all bandit problems $(\mathbb{P}_x)_{x \in \mathcal{Q}_1}$
 with distributions \mathbb{P}_x
 having support \mathcal{Q}_1 .

Indeed: Consider $(\delta_x)_{x \in \mathcal{Q}_1}$ the bandit problem in which each arm x is associated with the Dirac mass on x .

Fix any strategy: it gets $y_t = 0$ a.s. and uses a sequence of (possibly) random choices I_t , $t \geq 1$.

Since probability distributions can only have at most countably many atoms,

$$\mathcal{Y} = \{x \in \mathcal{Q}_1 : \exists t \mid \mathbb{P}\{I_t = x\} > 0 \text{ under } (\delta_x)_{x \in \mathcal{Q}_1}\}$$

is countable. In particular, $\mathcal{Q}_1 \setminus \mathcal{Y}$ is non-empty.

But the strategy behaves the same under the problem

$$(\tilde{\nu}_x^1)_{x \in \mathcal{Q}_1} \quad \text{in which} \quad \begin{cases} \tilde{\nu}_x^1 = \delta_0 & \forall x \neq x_0 \\ \tilde{\nu}_{x_0}^1 = \delta_1 & \text{for one fixed } x_0 \in \mathcal{Q}_1 \setminus \mathcal{Y} \end{cases}$$

With probability 1, it thus never hits x_0 .

$$\text{Therefore, } y_t = 0 \text{ a.s. a.s. and } \bar{R}_T = \frac{T}{T} - \mathbb{E}\left[\sum_{t=1}^T y_t\right] = 0.$$

Actually, continuity is

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

Theorem: metric space and let \mathcal{F} be the set of bandit problems $(\mathbb{P}_x)_{x \in A}$

with: $\rightarrow \forall x, \mathbb{P}_x$ is a distribution over \mathcal{Q}_1

\rightarrow a continuous mean-payoff function $f: x \mapsto E(\mathbb{P}_x)$

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

Corollary. Let \mathcal{F}^{all} be the family of all bandit models $(\mathbb{P}_x)_{x \in A}$ with distributions \mathbb{P}_x over $[0,1]$. Then the regret against \mathcal{F}^{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}^{lip} be the family of bandit models $(\mathbb{P}_x)_{x \in [0,1]}$ with distributions \mathbb{P}_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 \nearrow a performance bound by splitting $[0,1]$ into $[0,1/K], [1/K, 2/K], \dots, [i-1/K, i/K]$ with $i=1, \dots, K$ for a fixed K , where each bin $[i-1/K, i/K]$ plays the role of an arm 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.

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

where L is the Lipschitz constant of the mean-payoff function of \mathbb{P}_x

Proof of the Corollary:

We endow A with the discrete topology, ie, choose the distance $d(x,y) = 1$ if $x \neq y$. 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 \mathcal{A} is separable:

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

We pick actions in a triangular fashion:

Regime 1: UCB based on x_1, x_2
(fresh start)

$I_1^{(1)} \dots I_k^{(1)}$

Regime r : UCB based on x_1, x_2, \dots, x_r
(fresh start)

$I_1^{(r)} \dots I_{(r+1)^2}^{(r)}$

In Regime r :

$$(r+1)^2 \max_{s \leq r} f(x_s) - E \left[\sum_{t=S_r+1}^{S_r + (r+1)^2} y_t \right] \quad (*)$$

Starts at time
 $S_r + 1 = 2^2 + \dots + r^2 + 1$

$$\leq C \sqrt{r^3 \ln r}$$

well-chosen numerical constant

order of magnitude of the distribution-free regret bound for UCB on $(r+1)^2$ steps with r arms
(we saw this bound as an exercise)

Now, let $\varepsilon > 0$ and let $\tilde{r}_\varepsilon \in \mathbb{N}^*$ s.t. $\mu_{x_{\tilde{r}_\varepsilon}} \geq \sup_A f - \varepsilon$

(\tilde{r}_ε exists by separability of \mathcal{A} and continuity of f)

In particular,

$$\max_{s \leq \tilde{r}_\varepsilon} f(x_s) \geq \sup_A f - \varepsilon \quad (**)$$

We denote by r_T the index of the regime where T lies:

we have that S_r is of the order of r^3

thus that r_T is of the order of $T^{1/3}$

in particular, $r_T = O(T^{1/3})$

The regret can be decomposed (for T large enough) as

$$\begin{aligned}
 \bar{R}_T &= T \sup_A f - E\left[\sum_{t=1}^T y_t\right] = \text{sum of the regrets of each regime} \\
 &= \underbrace{\sum_{r=1}^{r_\Sigma-1} (r\epsilon)^2}_{\substack{\text{initial regimes,} \\ \text{regret bounded} \\ \text{by their lengths}}} + \underbrace{\sum_{r=r_\Sigma}^{r_T-1} \left((r\epsilon)^2 \Sigma + c\sqrt{r^3 \ln r} \right)}_{\substack{\text{if bounds } (*) \\ \text{and } (**)}} + \underbrace{(r_T+1)^2}_{\substack{\text{regime } r_T \\ \text{may be} \\ \text{incomplete}}} \\
 &= O(1) \\
 &\leq TE + c \sum_{r=r_\Sigma}^{r_T-1} r^{3/2} \sqrt{\ln r} = O(T^{2/3}) \\
 &\leq TE + O(r_T^{5/2} \sqrt{\ln r_T}) \\
 &= TE + O(T^{5/6} \sqrt{\ln T})
 \end{aligned}$$

All in all, $\limsup \frac{\bar{R}_T}{T} \leq \varepsilon$, which is true $\forall \varepsilon > 0$,

that is, $\lim \frac{\bar{R}_T}{T} = 0$ 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.

$$\rho = \inf \{ d(x_1, x_2) : x_1, x_2 \in \mathcal{D} \} > 0.$$

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

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

* We consider the bandit models $\mathbb{J}^{(a)}$ inducing mean-payoff functions

$f^{(a)} : x \in A \mapsto (1 - \frac{d(x, a)}{\rho})^+$; in particular, $\mathbb{J}^{(a)}_x = s_0$ for $x \notin B(a, \rho)$.

We proceed as in the example showing the necessity of continuity when $A = (\mathbb{Q})$

and consider the bandit model $(s_0)_{x \in A}$ as well as any strategy

and the laws induced by the I_t under this model: let d_t

be the law of I_t under $(s_0)_{x \in A}$ and let $d = \sum_{t \geq 1} \frac{1}{2^t} d_t$.

[as only countably many balls can have a positive mass under d]

\hookrightarrow There exists $a \in A$ s.t. $d(B(a, \rho)) = 0$, that is, s.t.,

$$\forall t \geq 1, \quad \mathbb{P}(I_t \in B(a, \rho) \text{ under } (s_0)_{x \in A}) = 0.$$

The considered strategy is therefore such that the I_t have

the same distribution under $(s_0)_{x \in A}$ and $\mathbb{J}^{(a)}$. In particular,

$E[\sum_{t=1}^T Y_t] = 0$ in both cases, but in the latter case,

$\sup f^{(a)} = 1$, so that $\bar{R}_T = T$ against $\mathbb{J}^{(a)}$. The regret

is not controlled against $\mathbb{J}^{(a)} \in \mathbb{J}^{\text{cont}}$

[Back to K-armed bandits]

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

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

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

One can show matching upper and lower bounds (with associated strategies):

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

where

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

We will only prove the lower bound part

but not discuss the strategy, called KL-UCB, to achieve the bound.

Kullback
Leibler
divergence

expectation
of $\tilde{\pi}_a^*$

} it's just a matter of time we would need about 3h to discuss this strategy

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

Part 3: The Kullback-Leibler divergence, definition and first properties

The Kullback-Leibler divergence: definition and basic properties.

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

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

Basic facts:

- Existence of the defining integral when $P \ll Q$: because $\Psi: x \mapsto x \ln x$ is bounded from below on $[0, +\infty)$
- $KL(P||Q) \geq 0$ and $0 = KL(P||Q)$ if and only if $P = Q$:

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

$$KL(P||Q) = \int_{\Omega} \Psi\left(\frac{dP}{dQ}\right) dQ \geq \Psi\left(\underbrace{\int \frac{dP}{dQ} dQ}_{=1}\right) = 0$$

with equality if and only

if $\frac{dP}{dQ}$ is Q -a.s. constant, i.e., $P = Q$

Exercise : A useful rewriting. Prove the following result:

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

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

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

Beware: with the usual measure-theoretic conventions, if $x \neq 0$ and $y = 0$, then $\frac{x}{y} \neq \frac{y}{x}$ ↳ 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 (Ω, \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

$$P^X \ll Q^X, \quad \text{with} \quad \frac{dP^X}{dQ^X} = E_Q \left[\frac{dP}{dQ} \mid X = \cdot \right] \stackrel{\text{def.}}{=} \gamma$$

$$\text{i.e., } \gamma(x) = E_Q \left[\frac{dP}{dQ} \mid X = x \right].$$

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

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

$$\stackrel{\text{def.}}{=} \int_{\Omega} \mathbb{1}_B(x) \gamma(x) dQ = \int_{\Omega} \mathbb{1}_B \gamma dQ.$$

Therefore,

$$\begin{aligned} KL(P^X, Q^X) &= \int_{\Omega} \gamma \ln \gamma dQ^X = \int_{\Omega} \gamma(x) \ln \gamma(x) dQ \\ &= \int_{\Omega} \left(E_Q \left[\frac{dP}{dQ} \mid X \right] \ln E_Q \left[\frac{dP}{dQ} \mid X \right] \right) dQ \quad \stackrel{\text{def.}}{\longrightarrow} \text{definition of } \gamma \\ &\leq \int_{\Omega} E_Q \left[\frac{dP}{dQ} \ln \frac{dP}{dQ} \mid X \right] dQ \quad \stackrel{\text{conditional version of Jensen's inequality}}{\longrightarrow} \\ &\stackrel{\text{tower rule}}{\longrightarrow} = \int_{\Omega} \left(\frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ = KL(P, Q) \end{aligned}$$

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; Bauceron, Ligos, 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 $(x,y) \in [0,+\infty)^2 \mapsto \left(\frac{x}{y} \ln \frac{x}{y}\right)$

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

Corollary (joint convexity of KL): For all probability distributions P_1, P_2 and Q_1, Q_2 over the same measurable space (Ω, \mathcal{F}) , and all $d \in (0,1)$,

$$\text{KL}\left((1-d)P_1 + dP_2, (1-d)Q_1 + dQ_2\right) \leq (1-d)\text{KL}(P_1, Q_1) + d\text{KL}(P_2, Q_2)$$

Proof: We augment (Ω, \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 \bar{P} be a probability measure on $(\Omega \times \{1,2\}, \mathcal{F}'')$ such that

$$\begin{cases} J \sim \text{Ber}(d) \\ X | J=j \sim P_j \end{cases} \quad \begin{matrix} \text{(and a similar definition for } Q \text{)} \\ \text{based on } Q_1, Q_2 \end{matrix}$$

that is, $\forall j \in \{1,2\} \quad \forall A \in \mathcal{F} \quad \bar{P}(A \times \{j\}) = ((1-d)1_{\{j=1\}} + d1_{\{j=2\}}) P_j(A)$

$$\text{Now, } P^x = (1-d)P_1 + dP_2$$

$$Q^x = (1-d)Q_1 + dQ_2$$

and (as we prove below) $\text{KL}(P, Q) = (1-d)\text{KL}(P_1, Q_1) + d\text{KL}(P_2, Q_2)$
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

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

This entails that

$$\begin{aligned} \text{KL}(P, Q) &= \int_{\Omega \times \{1,2\}} \left(\frac{dP}{dQ}(w, j) \ln \frac{dP}{dQ}(w, j) \right) dQ(w, j) \\ &\quad \text{(we just use that for } f \geq 0 \text{ constant)} \\ &= \int_{\Omega \times \{1,2\}} \left(\frac{dP}{dQ}(w, 1) \ln \frac{dP}{dQ}(w, 1) \right) 1_{\Omega \times \{1\}}(w, j) dQ(w, j) \\ &\quad + \int_{\Omega \times \{1,2\}} \left(\frac{dP}{dQ}(w, 2) \ln \frac{dP}{dQ}(w, 2) \right) 1_{\Omega \times \{2\}}(w, j) dQ(w, j) \\ &= \underbrace{\int_{\Omega} \left(\frac{dP_1}{dQ_1}(w) \ln \frac{dP_1}{dQ_1}(w) \right) (1-d)dQ_1(w)}_{= \text{KL}(P_1, Q_1) \times (1-d)} + \underbrace{\int_{\Omega} \dots}_{d\text{KL}(P_2, Q_2)} \end{aligned}$$

$\int f d\mu = \int f 1_A d\mu$
+ if the dQ is integrable or not

KL for product measures. (\Leftrightarrow The independent case)

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

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