

Part 1: The Kullback-Leibler divergence, more properties

Recap of what we already saw on KL-divergences:

Definition:
$$KL(P, Q) = \begin{cases} +\infty & \text{if } P \not\ll Q \\ \int_{\Omega} \left(\frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ & \text{if } P \ll Q \end{cases}$$

First properties: $KL(P, Q) \geq 0$ with $KL(P, Q) = 0 \iff P = Q$

Data-processing inequality (with X r.v.): $KL(P^X, Q^X) \leq KL(P, Q)$

Joint convexity of KL:
$$KL((1-d)P_1 + dP_2, (1-d)Q_1 + dQ_2) \leq (1-d)KL(P_1, Q_1) + dKL(P_2, Q_2)$$

First new result today:

KL for product measures

(\leftrightarrow the independent case, while the dependent case will be considered in a few pages)

Prop: 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:

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

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 $\int (f \cdot g) d\mu = \int f d\mu \int g d\mu$

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

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

$$= KL(P, Q)$$

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 (Ω, \mathcal{F})

Let $X: (\Omega, \mathcal{F}) \rightarrow ([a, 1], \mathcal{B}([a, 1]))$ be any $[a, 1]$ -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 $[a, 1]$ and augment the underlying measurable space into $(\Omega \times [a, 1], \mathcal{F} \otimes \mathcal{B}([a, 1]))$, over which we consider the product-distributions $P \otimes m$ and $Q \otimes m$.

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

$$\begin{aligned}
 \underbrace{KL\left(\underbrace{(P \times m)}^{\mathbb{1}_E}, \underbrace{(Q \times m)}^{\mathbb{1}_E}\right)}_{\text{Ber}(P \times m(E)) \quad \text{Ber}(Q \times m(E))} &\leq KL(P \times m, Q \times m) \\
 &= KL(P, Q) + KL(m, m) \\
 &\stackrel{\substack{\uparrow \\ \text{if product} \\ \text{distributions}}}{=} KL(P, Q)
 \end{aligned}$$

Thus: $KL(\text{Ber}(P \times m(E)), \text{Ber}(Q \times m(E))) \leq KL(P, Q)$

The proof is concluded by picking $E \in \mathcal{F} \otimes \mathcal{B}([0,1])$ such that $P \times m(E) = \mathbb{E}_P[X]$ and $Q \times m(E) = \mathbb{E}_Q[X]$

Namely, $E = \{(\omega, x) \in \Omega \times [0,1] : x \leq X(\omega)\}$

By Tonelli's theorem:

$$\begin{aligned}
 P \times m(E) &= \int_{\Omega} \left(\int_{[0,1]} \mathbb{1}_{\{x \leq X(\omega)\}} dm(x) \right) dP(\omega) \\
 &= \int_{\Omega} X(\omega) dP(\omega) = \mathbb{E}_P[X]
 \end{aligned}$$

and a similar equality for $Q \times m(E)$.

The chain rule — A generalization of the decomposition of the KL between product-distributions.

We will need it in a special case only, when the joint distributions follows from one of the marginal distributions via a stochastic kernel.

Definition: Let (Ω, \mathcal{F}) and (Ω', \mathcal{F}') be two measurable spaces; we denote by $\mathcal{P}(\Omega', \mathcal{F}')$ the set of probability measures over (Ω', \mathcal{F}') .

A stochastic kernel K is a mapping $(\Omega, \mathcal{F}) \rightarrow \mathcal{P}(\Omega', \mathcal{F}')$
(regular) $\omega \mapsto K(\omega, \cdot)$

such that $\forall B \in \mathcal{F}' \quad \omega \mapsto K(\omega, B)$ is \mathcal{F} -measurable.

Now, consider two such kernels K and L , and two probability measures P and Q over (Ω, \mathcal{F}) . Then KP and LQ defined below are probability measures over $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$, by some extension theorem: (Carathéodory?)

$\forall A \in \mathcal{F}, \forall B \in \mathcal{F}'$

$$KP(A \times B) = \int_{\Omega} \mathbb{1}_A(\omega) \underbrace{K(\omega, B)}_{\text{is indeed measurable}} dP(\omega)$$

$$LQ(A \times B) = \int_{\Omega} \mathbb{1}_A(\omega) L(\omega, B) dQ(\omega)$$

An extension of Fubini (-Tonelli) Theorem

immediate extension to φ lower bounded by $\mu \in \mathbb{R}$

Lemma: Let $\varphi: \Omega \times \Omega' \rightarrow \mathbb{R}$ be $\mathcal{F} \otimes \mathcal{F}'$ -measurable and either ≥ 0 or KP -integrable.

Then $\omega \mapsto \int_{\Omega'} \varphi(\omega, \omega') K(\omega, d\omega')$ is \mathcal{F} -measurable and

$$\int_{\Omega \times \Omega'} \varphi dKP = \int_{\Omega} \left(\int_{\Omega'} \varphi(\omega, \omega') K(\omega, d\omega') \right) dP(\omega)$$

Proof: (sketch) The result is true for $\varphi = \mathbb{1}_{A \times B}$ by definition of KP including measurability of $\omega \mapsto \int \varphi(\omega, \cdot) K(\omega, d\cdot)$ by regularity of K .
 Extension to $\mathbb{1}_E$ for any $E \in \mathcal{F} \otimes \mathcal{F}'$ by an argument of monotone convergence (including the $\omega \mapsto \int_{\Omega'} \dots$ measurability) or algebra contained / monotone class theorem, using.
 Extension to $\varphi \geq 0$ by monotone convergence.
 $\varphi \in L^1$

Question: Does anyone have a simpler argument?

Theorem [chain rule for KL]:

As soon as (*) $K(\omega, \cdot) \ll L(\omega, \cdot)$ for P -almost all $\omega \in \Omega$

with (**) the existence of a version $g: (\omega, \omega') \mapsto \frac{dK(\omega, \cdot)}{dL(\omega, \cdot)}(\omega')$ being $\mathcal{F} \otimes \mathcal{F}'$ -measurable,

Then

$$KL(KP, LQ) = KL(P, Q) + \int_{\Omega} KL(K(\omega, \cdot), L(\omega, \cdot)) dP(\omega)$$

where $\omega \mapsto KL(K(\omega, \cdot), L(\omega, \cdot))$ is indeed \mathcal{F} -measurable and ≥ 0 so that the integral in the right-hand side is well defined.

Remark: see a remark stated in two pages for the (lack of) necessity of Assumptions (*) and (**).

Proof: * By bi-measurability of $g \ln g$, and since $g \ln g$ is lower bounded, (an immediate extension of) the previous lemma can be applied to get

$$\omega \mapsto \int_{\Omega'} g(\omega, \cdot) \ln(g(\omega, \cdot)) L(\omega, d\cdot) = KL(K(\omega, \cdot), L(\omega, \cdot))$$

is \mathcal{F} -measurable and ≥ 0 , with:

we will not use this, actually

$$\int_{\Omega \times \Omega'} g \ln g \, dLIP = \int KL(K(\omega, \cdot), L(\omega, \cdot)) \, dP(\omega)$$

* We assume $P \ll Q$, let $f = \frac{dP}{dQ}$: what can we say about $(\omega, \omega') \mapsto f(\omega) g(\omega, \omega')$?

$$\begin{aligned} & \int \mathbb{1}_{A \times B}(\omega, \omega') f(\omega) g(\omega, \omega') \, dLQ(\omega, \omega') \\ & \stackrel{\text{of extension of Tonelli}}{=} \int_{\Omega} \left(\int_{\Omega'} \mathbb{1}_B(\omega') g(\omega, \omega') L(\omega, d\omega') \right) \mathbb{1}_A(\omega) f(\omega) \, dQ(\omega) \\ & \qquad \qquad \qquad = \int_{\Omega'} \mathbb{1}_B(\omega') K(\omega, d\omega') \\ & \qquad \qquad \qquad = K(\omega, B) \quad \text{given the definition of } g \\ & = \int \underbrace{\mathbb{1}_A(\omega) K(\omega, B)}_{\mathcal{F}\text{-measurable}} \underbrace{f(\omega) \, dQ(\omega)}_{\text{since } f = \frac{dP}{dQ}} = KP(A \times B) \quad \text{by def. of } KP \end{aligned}$$

By Radon-Nikodym's Theorem:

$$\frac{dKP}{dLQ} = fg \quad LQ\text{-as}$$

* It is easily seen that $KP \ll LQ \Rightarrow P \ll Q$ (in all cases, even without (*) and (**))

* Therefore, we have $KP \ll LQ \Leftrightarrow P \ll Q$ under (*) and (**), we thus assume with no loss of generality that $KP \ll LQ$ and $P \ll Q$ (otherwise, both $= +\infty$ and the putative equality is $+\infty = +\infty$).

Then,

$$KL(KP, LQ) = \int_{\Omega \times \Omega'} (f(\omega) g(\omega, \omega') \ln f(\omega) g(\omega, \omega')) \, dLQ(\omega, \omega')$$

what is being integrated in the \int_{Ω} integral is a function bounded on Ω'

The lemma (extension of Fubini-Tonelli) extends to it:

$$\int_{\Omega} (fg \ln(fg)) d\mathbb{Q} = \int_{\Omega} f(\omega) \left(\int_{\Omega'} (g(\omega, \omega') \ln g(\omega, \omega') + g(\omega, \omega') \ln f(\omega)) L(\omega, d\omega') \right) d\mathbb{Q}(\omega)$$

$$= \int_{\Omega} \left(\int_{\Omega'} (g(\omega, \omega') \ln g(\omega, \omega')) L(\omega, d\omega') + \underbrace{(\ln f(\omega))}_{=1} g(\omega, \omega') L(\omega, d\omega') \right) f(\omega) d\mathbb{Q}(\omega)$$

$KL(K(\omega, \cdot), L(\omega, \cdot))$

$$= \int_{\Omega} \left(KL(K(\omega, \cdot), L(\omega, \cdot)) + \ln f(\omega) \right) f(\omega) d\mathbb{Q}(\omega)$$

$$= \int_{\Omega} \underbrace{KL(K(\omega, \cdot), L(\omega, \cdot)) f(\omega) d\mathbb{Q}(\omega)}_{d\mathbb{P}(\omega)} + \int_{\Omega} \underbrace{f(\omega) \ln f(\omega) d\mathbb{Q}(\omega)}_{= KL(\mathbb{P}, \mathbb{Q})}$$

Sum of two functions bounded below

REMARKS ON THE ASSUMPTIONS

- The assumptions (*) and (**) will be satisfied for the applications we have in mind
- They can be relaxed: it suffices to assume that Ω' is a topological space with a countable base (a "second-countable space") and \mathcal{F} is the Borel σ -algebra.

I.e., there exists some countable collection $(O_n)_{n \geq 1}$ of open sets of Ω' such that each open set V of Ω' can be written

$$V = \bigcup_{i: O_i \in V} O_i$$

that is, as a countable union of elements of $(O_n)_{n \geq 1}$.

Ex: Ω' a separable metric space \rightarrow we will consider $\Omega' = [0,1] \times (\mathbb{R} \times [0,1])^{\mathbb{N}}$

\hookrightarrow See details in the additional document.

CREDITS: Marin Billu + Hadi Hadjici, M2 students of Spring 2017

Part 2: Regret lower bound for stochastic bandits

Lower bounds on the regret for stochastic bandits.

Here is first a summary of the setting and context of stochastic bandits:

- K arms each indexed by $a = 1, 2, \dots, K$
- With each arm is associated a probability distribution $\nu_a \in \mathcal{D}$
- \mathcal{D} is the bandit model: a subset of $\mathcal{M}_1(\mathbb{R})$, the set of probability distributions over \mathbb{R} with an expectation
- A bandit problem is denoted by $\nu = (\nu_a)_{a \in \{1, \dots, K\}}$
- Important quantities and relation:

$\mu_a = E(\nu_a)$ is the expectation of ν_a

$\mu^* = \max_{a=1, \dots, K} \mu_a$ is the largest expectation within ν

$\Delta_a = \mu^* - \mu_a$ is the gap for arm a

Arm a is suboptimal if $\Delta_a > 0$

$U_0, U_1, U_2, \dots, U_{t-1}$
i.i.d. $\sim U_{[0,1]}$

- Protocol: at each round $t = 1, 2, \dots$

1. The decision-maker picks $I_t \in \{1, \dots, K\}$ possibly at random based on an auxiliary randomization U_{t-1}
2. She gets a reward Y_t drawn at random according to $\nu_{I_t}^*$ (given I_t); this is the only piece of information she gets.

(Note: The decision-maker knows the specific ν_a at hand)

- Aim/regret: maximize $E\left[\sum_{t=1}^T Y_t\right]$
which is equivalent to minimizing (controlling from above)
 $R_T = T\mu^* - E\left[\sum_{t=1}^T Y_t\right]$

- Rewriting by tower rule:

$$R_T = T\mu^* - E\left[\sum_{t=1}^T \mu_{I_t}\right] = \sum_{a=1}^K \Delta_a E[N_a(T)]$$

where $N_a(T) = \sum_{t=1}^T \mathbb{1}_{I_t=a}$ is the number of times arm a was pulled between 1 and T

! It is thus necessary and sufficient to control $E[N_a(T)]$ for suboptimal arms a

- What is a (randomized) strategy?

A sequence of measurable functions $(\Psi_t)_{t \geq 0}$ with

$$\Psi_t: \underbrace{H_t = (U_0, Y_1, U_1, \dots, Y_t, U_t)}_{\text{history for the first } t \text{ rounds}} \rightarrow \underbrace{\Psi_t(H_t) = I_{t+1}}_{\text{arm picked at round } t+1}$$

- Strategies that are consistent w.r.t. a model \mathcal{D} :

If for all bandit problems $\vec{\nu} \in \mathcal{D}^k$,

$$\forall \epsilon \in (0, 1], \quad \forall \Delta_n \text{ s.t. } \Delta_n \geq 0, \quad \mathbb{E}[N_\epsilon(T)] = o(T^\Delta).$$

- Result: For "well-behaved" models \mathcal{D} , there exist consistent strategies.

E.g.: at least $\mathcal{D} = \mathcal{M}_1([0, 1])$, see the UCB strategy.

- Typical bounds for good strategies (stated in an asymptotic way, even though non-asymptotic bounds are available)

$\forall \vec{\nu} \in \mathcal{D}^k, \quad \forall \Delta_n \text{ s.t. } \Delta_n \geq 0,$

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}[N_\epsilon(T)]}{\ln T} \leq C_\Delta(\vec{\nu})$$

where $C_\Delta(\vec{\nu})$ is a problem-dependent constant.

- Optimal (in some sense) such constant: $C_\Delta(\vec{\nu}) = \frac{1}{\text{Kinf}(\vec{\nu}_a, \mu^*; \mathcal{D})} = \frac{1}{\text{Kinf}(\vec{\nu}_a, \mu^*)}$

where $\text{Kinf}(\vec{\nu}_a, \mu^*; \mathcal{D}) = \text{Kinf}(\vec{\nu}_a, \mu^*) = \inf \left\{ \text{KL}(\vec{\nu}_a, \vec{\nu}_a') : \begin{array}{l} \vec{\nu}_a' \in \mathcal{D} \\ \mathbb{E}[\vec{\nu}_a'] \geq \mu^* \end{array} \right\}$

with the convention: $\inf \emptyset = +\infty$.

We will only prove one part of this optimality: a lower bound on $C_\Delta(\vec{\nu})$.

Theorem:

For all bandit models $\mathcal{D} \subset \mathcal{M}_1(\mathbb{R})$,

(see Lai and Robbins, 1985; Burnetas and Katehakis, 1996)

For all strategies Ψ consistent w.r.t. \mathcal{D} (possibly randomized),

For all bandit problems $\vec{\nu} = (\nu_a)_{a \in \{1, \dots, k\}} \in \mathcal{D}^k$,

For all suboptimal arms a (ie, such that $\Delta_n > 0$),

$$\liminf \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\text{Kinf}(\nu_a, \mu^*; \mathcal{D})}$$

Corollary: For all bandit models $\mathcal{D} \subseteq \mathcal{U}_1(\mathbb{R})$,
 For all (possibly randomized) strategies ψ consistent w.r.t \mathcal{D} ,
 For all bandit problems $\vec{\nu} = (\nu_a)_{a \in \{1, \dots, K\}} \in \mathcal{D}^K$,

$$\liminf_{T \rightarrow \infty} \frac{\bar{R}_T}{\ln T} \geq \sum_{a: \Delta_a > 0} \frac{\Delta_a}{K_{\text{inf}}(\vec{\nu}_a, \mu^*, \mathcal{D})}.$$

To prove this theorem (and to prove other lower bounds), we will need the following fundamental inequality. In its statement, $\mathbb{P}_\psi^{\vec{\nu}}$ and $\mathbb{E}_\psi^{\vec{\nu}}$ refer to the probability distribution and the expectation induced by the bandit problem $\vec{\nu} \in \mathcal{D}^K$.

Example: $\mathbb{P}_\psi^{\vec{\nu}}$ is the law of $H_T = (U_0, Y_1, U_1, \dots, Y_T, U_T)$ when the bandit problem is $\vec{\nu}$. Actually, $\mathbb{P}_\psi^{\vec{\nu}}$ strongly depends on the strategy ψ used but we omit this dependency in the notation.

Lemma (Fundamental inequality for stochastic bandits):

For all bandit problems $\vec{\nu} = (\nu_a)_{a \in \{1, \dots, K\}}$ and $\vec{\nu}' = (\nu'_a)_{a \in \{1, \dots, K\}}$ in \mathcal{D}^K
 with $\vec{\nu}' \ll \vec{\nu}$ for all a ,

For all strategies
 For all random variables Z taking values in $[0, 1]$ and that are $\sigma(H_T)$ -measurable,

$$\sum_{a=1}^K \mathbb{E}_\psi^{\vec{\nu}}[N_a(T)] \text{KL}(\nu_a, \nu'_a) = \text{KL}(\mathbb{P}_\psi^{\vec{\nu}}, \mathbb{P}_\psi^{\vec{\nu}'}) \geq \text{KL}(\text{Ber}(\mathbb{E}_\psi^{\vec{\nu}}[Z]), \text{Ber}(\mathbb{E}_\psi^{\vec{\nu}'}[Z]))$$

⚠ The dependence on the strategy is hidden in the $\mathbb{E}_\psi^{\vec{\nu}}[N_a(T)]$, $\mathbb{E}_\psi^{\vec{\nu}}[Z]$ and $\mathbb{E}_\psi^{\vec{\nu}'}[Z]$

NOTE: This lemma is our key to perform an implicit change of measures in the proof of the theorem.

Proof of the theorem (based on the lemma) We have $K_{\text{inf}}(\nu_a, \mu^*) = \inf \{ KL(\nu_a, \nu_a^i) : \nu_a^i \in \mathcal{D} \text{ and } E(\nu_a^i) \geq \mu^* \}$
 $= \inf \{ KL(\nu_a, \nu_a^i) : \nu_a^i \in \mathcal{D}, \nu_a \ll \nu_a^i \text{ and } E(\nu_a^i) \geq \mu^* \}$
 (cf. convention: $\inf \emptyset = +\infty$ and the fact that $KL(\nu_a, \nu_a^i) = +\infty$ when $\nu_a \not\ll \nu_a^i$)

This is why we will

- Fix $\mathcal{D}, \Psi, \bar{\nu}$ and a s.t. $\Delta_a > 0$
- Fix an alternative model ν^3 of the form

$$\begin{cases} \nu_k^3 = \bar{\nu}_k & \forall k \neq a \\ \nu_a^3 & \text{s.t. } \nu_a^3 \in \mathcal{D}, \bar{\nu}_a \ll \nu_a^3 \text{ and } E(\nu_a^3) \geq \mu^* \end{cases}$$

That is, $\bar{\nu}$ and ν^3 only differ at a ; a is the unique optimal arm in ν^3

- Take $Z = N_a(T)/T$ which is indeed $[0,1]$ -valued $\sigma(H_T)$ -measurable.

Our fundamental inequality yields, since $\bar{\nu}$ and ν^3 only differ at a :

$$E_{\bar{\nu}}[N_a(T)] KL(\bar{\nu}_a, \nu_a^3) \geq KL(\text{Ber}(E_{\bar{\nu}}[N_a(T)/T]), \text{Ber}(E_{\nu^3}[N_a(T)/T])) \\ \geq -\ln 2 + (1 - E_{\bar{\nu}}[N_a(T)/T]) \ln \frac{1}{1 - E_{\nu^3}[N_a(T)/T]}$$

indeed: $KL(\text{Ber}(p), \text{Ber}(q))$

$$= p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q}$$

$$= \underbrace{p \ln \frac{1}{q}}_{\geq 0} + (1-p) \ln \frac{1}{1-q} + \underbrace{(p \ln p + (1-p) \ln(1-p))}_{\geq -\ln 2 \text{ by a simple function study over } [0,1]}$$

$$\geq -\ln 2 + (1-p) \ln \frac{1}{1-q}$$

for all $p, q \in (0,1)$ and even for all $p, q \in [0,1]$ (study the cases $q=0$ and $q=1$ separately)

Now, the considered strategy Ψ is consistent and:

- in the problem ν^3 , a is suboptimal: $E_{\nu^3}[N_a(T)/T] \rightarrow 0$

— in the problem \mathcal{J}' , all arms $k \neq a$ are suboptimal:

$$\text{for all } \alpha \in (0, 1], \quad T - \mathbb{E}_{\mathcal{J}'}[N_a(T)] = \sum_{k \neq a} \mathbb{E}_{\mathcal{J}'}[N_k(T)] = o(T^\alpha)$$

↳ in particular, for T large enough,

$$\frac{1}{1 - \mathbb{E}_{\mathcal{J}'}[N_k(T)/T]} = \frac{T}{T - \mathbb{E}_{\mathcal{J}'}[N_k(T)]} \geq \frac{T}{T^\alpha} = T^{1-\alpha}$$

Substituting back and dividing by $\ln T$: for all $\alpha \in (0, 1]$, for T large enough:

$$\frac{\mathbb{E}_{\mathcal{J}'}[N_k(T)]}{\ln T} \text{KL}(\mathcal{J}'_{a_1}, \mathcal{J}'_{a_2}) \geq \frac{-\ln 2}{\ln T} + \underbrace{\left(1 - \mathbb{E}_{\mathcal{J}'}\left[\frac{N_k(T)}{T}\right]\right)}_{\rightarrow 0} \underbrace{\frac{\ln T^{1-\alpha}}{\ln T}}_{= 1-\alpha}$$

thus

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{J}'}[N_k(T)]}{\ln T} \text{KL}(\mathcal{J}'_{a_1}, \mathcal{J}'_{a_2}) \geq 1-\alpha$$

Letting $\alpha \rightarrow 0$:

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{J}'}[N_k(T)]}{\ln T} \text{KL}(\mathcal{J}'_{a_1}, \mathcal{J}'_{a_2}) \geq 1$$

Whether $\text{KL}(\mathcal{J}'_{a_1}, \mathcal{J}'_{a_2}) < +\infty$ or $= +\infty$, we thus get

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{J}'}[N_k(T)]}{\ln T} \geq \frac{1}{\text{KL}(\mathcal{J}'_{a_1}, \mathcal{J}'_{a_2})}$$

The left-hand side is independent of $\mathcal{J}'_{a_2} \in \mathcal{D}$ s.t. $\mathcal{J}'_{a_2} \succ \mathcal{J}'_{a_1}$ and $\mathbb{E}(\mathcal{J}'_{a_2}) \succ \mu^*$, so that taking the supremum of the right-hand side over these \mathcal{J}'_{a_2} , we get the desired $\frac{1}{\text{KL}(\mathcal{J}'_{a_1}, \mu^*)}$ lower bound.

Proof of the lemma:

- The inequality \geq is a direct application of the data-processing inequality w.r.t. expectations

- For the equality:

and similarly $P_{y_t}^{H_t} = K_T^*(\dots (K_1^* m) \dots)$

(i) We show by induction that $P_{y_t}^{H_t} = K_T(K_{T-1}(\dots (K_1 m) \dots))$

where K_t is the transition kernel:

$h \in [Q] \times (\mathbb{R} \times [Q])^{t-1} \mapsto K_t(h, \cdot) = \sum_{\psi_t(h)} \nu_a \otimes m$

we check below that it is regular

Indeed: $T=0: H_0 = U_0 \sim U_{[Q]}: P_{y_0}^{U_0} = m$

$t \rightarrow t+1: \forall A \in \mathcal{B}([Q] \times (\mathbb{R} \times [Q])^t) \quad \forall B' \in \mathcal{B}(\mathbb{R}), \forall B \in \mathcal{B}([Q]),$

$P_{y_t}^{H_t}(A \times B' \times B) = P_{y_t}(H_t \in A \text{ and } Y_{t+1} \in B' \text{ and } U_{t+1} \in B)$
 $= E_{y_t} [\mathbb{1}_A(H_t) P_{y_t}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B | H_t)]$

tower rule

definition of the branch of the strategy

$E_{y_t} [\mathbb{1}_A(H_t) \sum_{\psi_t(H_t)} \nu_a(B') m(B)]$

definition of K_{t+1}

$= E_{y_t} [\mathbb{1}_A(H_t) K_{t+1}(H_t, B' \times B)]$

$= \int \mathbb{1}_A K_{t+1}(h, B' \times B) dP_{y_t}^{H_t}(h)$

mere rewriting

$= K_{t+1} P_{y_t}^{H_t}(A \times B' \times B)$

definition of $K_{t+1} P_{y_t}^{H_t}$

- (2) We first check that the assumptions of the chain rule are satisfied:

* The K_t are regular transition kernels: $\forall E \in \mathcal{B}(\mathbb{R}) \otimes \mathcal{B}([Q]),$

$h \mapsto K_t(h, E) = \sum_{a=1}^K \mathbb{1}_{\{\psi_t(h)=a\}} \nu_a \otimes m(E)$

is measurable as ψ_t is measurable

* Assumption (*): $\forall h, K_t(h, \cdot) \ll K_t^*(h, \cdot)$ as $\forall a, \nu_a \ll \nu_a^*$ by assumption

* Assumption (**): $(h, (y,u)) \mapsto \frac{dK_t(h, \cdot)}{dK_t^*(h, \cdot)}(y,u)$
 is indeed \mathcal{F}_t -measurable
 (product of measurable functions) $= \sum_{a=1}^K \mathbb{1}_{\{Y_t(h)=a\}} \frac{d\gamma_a^*(y)}{d\gamma_a(y)}$

(3) We then may apply the chain rule and show by induction the desired result based on:

$$- \quad \text{KL}(\mathbb{P}_{\mathcal{Y}^t}, \mathbb{P}_{\mathcal{Y}^t}) = \text{KL}(\eta, \eta) = 0$$

$$\begin{aligned}
 - \quad \text{For } t \geq 0, \quad & \text{KL}(\mathbb{P}_{\mathcal{Y}^{t+1}}, \mathbb{P}_{\mathcal{Y}^{t+1}}) \\
 &= \text{KL}(\mathbb{K}_{t+1} \mathbb{P}_{\mathcal{Y}^t}, \mathbb{K}'_{t+1} \mathbb{P}_{\mathcal{Y}^t}) \\
 &= \text{KL}(\mathbb{P}_{\mathcal{Y}^t}, \mathbb{P}_{\mathcal{Y}^t}) + \int \text{KL}(\mathbb{K}_{t+1}(h, \cdot), \mathbb{K}'_{t+1}(h, \cdot)) d\mathbb{P}_{\mathcal{Y}^t}(h) \\
 &= \text{KL}(\mathbb{P}_{\mathcal{Y}^t}, \mathbb{P}_{\mathcal{Y}^t}) + \int \text{KL}(\gamma_{\mathcal{Y}^t}^*(h) \otimes \eta, \gamma_{\mathcal{Y}^t}^*(h) \otimes \eta) d\mathbb{P}_{\mathcal{Y}^t}(h) \\
 &= \text{KL}(\mathbb{P}_{\mathcal{Y}^t}, \mathbb{P}_{\mathcal{Y}^t}) + \sum_{a=1}^K \text{KL}(\gamma_a^*, \gamma_a^*) \int \mathbb{1}_{\{Y_t(h)=a\}} d\mathbb{P}_{\mathcal{Y}^t}(h) \\
 & \quad \underbrace{E[\mathbb{1}_{\{Y_t(H_t)=a\}}]} \\
 &= E[\mathbb{1}_{\{Y_{t+1}=a\}}]
 \end{aligned}$$

Exercise: $\frac{1}{K_{\text{inf}}(\mathcal{Z}_T, \mu^*, \mathcal{D})}$ vs. $\frac{8}{\Delta_n^2}$ for UCB

Recall that in the model $\mathcal{D} = \mathcal{J}(\mathcal{Q}, \mathcal{I})$, the UCB algorithm enjoys the following performance bound:

$$\forall i \in \mathcal{J}(\mathcal{Q}, \mathcal{I})^K, \quad \forall a \text{ s.t. } \Delta_n > 0, \\ \mathbb{E}_\mu [N_i(T)] \leq \frac{8}{\Delta_n^2} \ln T + 2.$$

Actually, there are refinements of UCB that get the distribution-dependent constant $\frac{8}{\Delta_n^2}$ arbitrarily close to $\frac{2}{\Delta_n^2}$.

But how do these $\frac{8}{\Delta_n^2}$ and $\frac{2}{\Delta_n^2}$ constants compare to $\frac{1}{K_{\text{inf}}(\mathcal{Z}_T, \mu^*, \mathcal{P}(\mathcal{Q}, \mathcal{I}))}$?

(1) For $p, q \in \mathcal{Q}, \mathcal{I}$, we denote

$$k(p, q) = \text{KL}(\text{Ber}(p), \text{Ber}(q))$$

Show that $\forall (p, q) \in \mathcal{Q}, \mathcal{I}^2, \quad k(p, q) \geq 2(p - q)^2$.

(2) Show Pinsker's inequality: let (Ω, \mathcal{F}) be a measurable space, let \mathbb{P}, \mathbb{Q} be two distributions over (Ω, \mathcal{F}) , then:

$$\| \mathbb{P} - \mathbb{Q} \|_{\text{TV}} = \sup_{A \in \mathcal{F}} | \mathbb{P}(A) - \mathbb{Q}(A) | \leq \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}, \mathbb{Q})}$$

↑
the total variation distance between \mathbb{P} and \mathbb{Q}

Even better, show the stronger form: $\sup_{Z: \mathcal{F}\text{-measurable taking values in } \mathcal{Q}, \mathcal{I}} | \mathbb{E}_{\mathbb{P}}[Z] - \mathbb{E}_{\mathbb{Q}}[Z] | \leq \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}, \mathbb{Q})}$

(3) Exhibit a lower bound on $K_{\text{inf}}(\mathcal{Z}_T, \mu^*, \mathcal{P}(\mathcal{Q}, \mathcal{I}))$ and conclude that some work is needed to get an upper bound matching our lower bound!

Exercise:Finite-time lower bound for small values of T

"All algorithms explore much!"

↳ We want to model that all algorithms must first explore uniformly all arms (\leftrightarrow exploration)

at least half of the time, before being able to perform exploitation more often.

(1) Establish the following local version of Pinsker's inequality:

$$- \forall 0 \leq p < q \leq 1, \quad \text{KL}(p, q) \geq \frac{1}{2 \max_{x \in [p, q]} x(1-x)} (p-q)^2 \geq \frac{1}{2q} (p-q)^2$$

- Why is it stronger than the global version of Pinsker's inequality?

(2) Show that all strategies smoothes than the uniform strategy [i.e. such that for all bandit problems, $\forall a$ s.t. $\mu_a = \mu^*$, $E[N_a(T)] \geq \frac{T}{K}$], we have:

$$\forall T \leq \frac{1}{8 K^*},$$

$$\text{where } K^* = \max_{j: \Delta_j > 0} K_{\text{inf}}(j, \mu_j^*)$$

$$\forall j \text{ s.t. } \Delta_j > 0,$$

$$E[N_j(T)] \geq \frac{1}{2} \frac{T}{K}$$

at least half of the time \uparrow uniform exploration

Hint: consider the same alternative bandit problems $\{j\}$ as in the theorem giving the asymptotic lower bound.