

Part 1: The Kullback-Leibler divergence, more properties

Recap of what we already saw on KL-divergences:

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

These properties heavily rely on the facts $\int_{\Omega} \frac{dP}{dQ} \geq 1$ and $\int_{\Omega} \frac{dP}{dQ} \ln \frac{dP}{dQ} \geq -\frac{1}{e}$
is strictly convex and

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

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

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

First new result today:

KL for product measures

(\iff 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:

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

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

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

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

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

we add $+2/\epsilon$

Therefore,
by Tonelli
+

we use that
if $f, g \geq 0$,
~~then~~
 $\int (f+g) d\mu$

$$= \int f d\mu + \int g d\mu$$

We ensure here
that $f, g \geq 0$ by the translations
by $\frac{1}{\epsilon}$

$$KL(P \ll P', Q \ll Q') = \int \left(\frac{dP}{dQ} \frac{dP'}{dQ'} \ln \left(\frac{dP}{dQ} \frac{dP'}{dQ'} \right) \right) d(Q \ll Q')$$

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

+ similar term
with $\ln \frac{dP'}{dQ}$, also
with $+2/\epsilon$

$$= KL(P', Q')$$

here
we apply Tonelli's
theorem (again because
 $\Omega \mapsto 2 \ln |\Omega|$ is lower bounded)

Consequence (Girvin, Nédélec, Stoltz, 2016):

Data-processing inequality with expectations of random variables

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

Let $X: (\Omega, \mathcal{F}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$ 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:

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

Proof: We denote by η 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 $P \otimes \eta$ and $Q \otimes \eta$.

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

$$\begin{aligned}
 \text{KL}\left(\underbrace{\text{P}(\alpha, m)}_{\text{Ber}(\text{P}(\alpha, m)(E))}, \underbrace{\text{Q}(\alpha, m)}_{\text{Ber}(\text{Q}(\alpha, m)(E))}\right) &\leq \text{KL}(\text{P}(\alpha, m), \text{Q}(\alpha, m)) \\
 &= \text{KL}(\text{P}, \text{Q}) + \text{KL}(m, m) \\
 &\stackrel{\uparrow}{=} \text{KL}(\text{P}, \text{Q})
 \end{aligned}$$

if product distributions

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

The proof is concluded by picking $E \in \mathcal{F} \times \mathcal{B}(\mathbb{R}^d)$ such that

$$\text{P}(\alpha, m)(E) = E_{\text{P}}[x] \quad \text{and} \quad \text{Q}(\alpha, m)(E) = E_{\text{Q}}[x]$$

Namely, $E = \{(w, x) \in \Omega \times \mathbb{R}^d : x \leq X(w)\}$

By Tonelli's theorem :

$$\begin{aligned}
 \text{P}(\alpha, m)(E) &= \int_{\Omega} \left(\int_{\mathbb{R}^d} \mathbb{1}_{\{x \leq X(w)\}} dm(x) \right) d\text{P}(w) \\
 &= \int_{\Omega} X(w) d\text{P}(w) = E_{\text{P}}[x]
 \end{aligned}$$

and a similar equality for $\text{Q}(\alpha, m)(E)$.

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) $w \mapsto K(w, \cdot)$

such that $\forall B \in \mathcal{F}'$ $w \mapsto K(w, B)$ 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 LP defined below are probability measures over $(\Omega \times \Omega', \mathcal{F} \times \mathcal{F}')$, by some extension theorem.
(Carathéodory?)

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

$$K\mathbb{P}(A \times B) = \int_{\Omega} \mathbf{1}_A(w) K(w, B) d\mathbb{P}(w)$$

is indeed measurable

$$L\mathbb{Q}(A \times B) = \int_{\Omega} \mathbf{1}_A(w) L(w, B) d\mathbb{Q}(w)$$

An extension of
Fubini (-Tonelli) Theorem
↓

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

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

$$\int_{\Omega \times \Omega'} \varphi dK\mathbb{P} = \int_{\Omega} \left(\int_{\Omega'} \varphi(w, w') K(w, dw') \right) d\mathbb{P}(w)$$

including measurability of $w \mapsto \int \varphi(w, \cdot) K(w, d\cdot)$ by regularity of K .

Proof: The result is true for $\varphi = \mathbf{1}_{A \times B}$ by definition of $K\mathbb{P}$

(sketch) Extension to $\mathbf{1}_E$ for any $E \in \mathcal{F} \otimes \mathcal{F}'$ by an argument of σ -algebra contained / monotone class theorem, using monotone convergence (including Fubini)

Extension to $\varphi \geq 0$ by monotone convergence
 $\varphi \in L^1$

$w \mapsto \int_{\Omega'} \dots$ measurability)

Question: Does anyone have a simpler argument?

actually with no loss of generality ✓

Theorem [chain rule for KL]: Assume $P \ll Q$

As soon as $(*) \quad K(w, \cdot) \ll L(w, \cdot)$ for Q almost all $w \in \Omega$

with $(**)$ the existence of a version $g: (\Omega, \mathcal{W}) \mapsto \frac{dK(w, \cdot)}{dL(w, \cdot)}(w)$
 being $\mathcal{F} \otimes \mathcal{F}'$ -measurable
 ✓ version up to a $L\mathbb{Q}$ -null set

Then

$$KL(K\mathbb{P}, L\mathbb{Q}) = KL(P, Q) + \int_{\Omega} KL(K(w, \cdot), L(w, \cdot)) d\mathbb{P}(w)$$

where $w \mapsto KL(K(w, \cdot), L(w, \cdot))$ 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

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

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

we will not use this, actually

$$\int_{\Omega \times \Omega'} g \ln g dLQ = \int_{\Omega} KL(K(\omega), L(\omega)) dLQ(\omega)$$

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

$$\begin{aligned} & \int \mathbb{1}_{A \times B}(\omega, \omega') f(\omega) g(\omega, \omega') dLQ(\omega, \omega') \\ &= \underbrace{\int_{\Omega} \left(\int_{\Omega'} \mathbb{1}_B(\omega') g(\omega, \omega') L(\omega, d\omega') \right) \mathbb{1}_A(\omega) f(\omega) dQ(\omega)}_{= \int_{\Omega'} \mathbb{1}_B(\omega') K(\omega, d\omega')} \\ &= K(A, B) \end{aligned}$$

given the definition of g

$$= \int \underbrace{\mathbb{1}_A(\omega) K(\omega, B)}_{\mathcal{F}\text{-measurable}} \underbrace{f(\omega) dQ(\omega)}_{\frac{dP(\omega)}{dQ(\omega)}} = KIP(A \times B) \quad \text{by def. of } KIP$$

By Radon-Nikodym's Theorem:

$$\frac{dKIP}{dLQ} = fg \quad LQ\text{-a.s}$$

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

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

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

$\varphi = fg \ln(fg)$ is lower bounded, Pre lemma (extension of Fubini-Tonelli) extends to it:

$$\int (\varphi g \ln(\varphi g)) dQ = \int f(\omega) \left(\int_{\Omega'} (g(\omega, \omega') \ln g(\omega, \omega') L(\omega, d\omega') + g(\omega, \omega') \ln f(\omega)) dQ(\omega') \right) dQ(\omega)$$

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

$$= \int \left(KL(K(\omega), L(\omega)) + \ln f(\omega) \right) f(\omega) dQ(\omega)$$

$$= \int \underbrace{KL(K(\omega), L(\omega))}_{dF(\omega)} f(\omega) dQ(\omega) + \int \underbrace{f(\omega) \ln f(\omega)}_{-KL(P, Q)} dQ(\omega)$$

Sum of two functions bounded from 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- σ -countable space") and \mathcal{F}' is the Borel σ -algebra.

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

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

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

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

↳ See details in the additional document.

CREDITS: Martin Brilu + Hedi Habliji, 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 $\pi_a \in \mathcal{D}$
- \mathcal{D} is the bandit model: a subset of $M_1(\mathbb{R})$, the set of probability distributions over \mathbb{R} with an expectation
- A bandit problem is denoted by $\pi = (\pi_a)_{a \in \{1, \dots, K\}}$
- Important quantities and notation:

$\mu_a = E(\pi_a)$ is the expectation of π_a

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

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

Arm a is suboptimal if $\Delta_a > 0$

U_0, U_1, U_2, \dots
iid $\sim U_{[0,1]}$



- 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 π_{I_t} (given I_t); this is the only piece of information she gets.

- Aim / regret:

$$\text{maximize } E\left[\sum_{t=1}^T y_t\right]$$

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 $N_a(T) = \sum_{t=1}^T \mathbf{1}_{\{I_t=a\}}$ 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

$$\Psi_t : H_t = (U_0, Y_1, U_1, \dots, Y_t, U_t) \mapsto \Psi_t(H_t) = I_{t+1}$$

history for the first
 t rounds

arm picked at
round $t+1$

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

If for all bandit problems $\mathcal{I} \in \mathbb{S}^K$,

$$\forall a \in \mathcal{I}, \quad \forall \alpha \text{ s.t. } \Delta_a > 0, \quad \mathbb{E}[N_a(T)] = o(T^\alpha).$$

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

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

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

$\forall \mathcal{I} \in \mathbb{S}^K, \quad \forall a \text{ s.t. } \Delta_a > 0,$

$$\limsup_{T \rightarrow +\infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \leq C_a(\mathcal{I})$$

where $C_a(\mathcal{I})$ is a problem-dependent constant.

- Optimal (in some sense) such constant: $C_a(\mathcal{I}) = \frac{1}{K_{\text{inf}}(\tilde{\nu}_a, \mu^*, \mathcal{D})} = \frac{1}{K_{\text{inf}}(\tilde{\nu}_a, \mu^*)}$

where $K_{\text{inf}}(\tilde{\nu}_a, \mu^*, \mathcal{D}) = K_{\text{inf}}(\tilde{\nu}_a, \mu^*) = \inf \left\{ \text{KL}(\tilde{\nu}_a, \nu'_a) : \begin{array}{l} \nu'_a \in \mathcal{D} \\ \mathbb{E}[\nu'_a] \geq \mu^* \end{array} \right\}$
with the convention: $\inf \emptyset = +\infty$.

We will only prove one part of this optimality: a lower bound on $C_a(\mathcal{I})$.

Theorem:

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

(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 $\mathcal{I} = (\nu_a)_{a \in \mathcal{I}, 1 \leq a \leq K} \in \mathbb{S}^K$,

For all suboptimal arms a (i.e., such that $\Delta_a > 0$),

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{K_{\text{inf}}(\tilde{\nu}_a, \mu^*, \mathcal{D})}$$

- 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 $\tilde{\gamma} = (\tilde{\gamma}_a)_{a \in \{1..K\}} \in \mathcal{D}^K$,

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

To prove this theorem (and to prove other lower bounds), we will need the following fundamental inequality. In its statement, P_T and E_T refer to the probability distribution and the expectation induced by the bandit problem $\tilde{\gamma} \in \mathcal{D}^K$.

Example: $P_T^{H_T}$ is the law of $H_T = (U_0, Y_1, U_1, \dots, Y_T, U_T)$ when the bandit problem is $\tilde{\gamma}$. Actually, $P_T^{H_T}$ 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 $\tilde{\gamma} = (\tilde{\gamma}_a)_{a \in \{1..K\}}$ and $\tilde{\gamma}' = (\tilde{\gamma}'_a)_{a \in \{1..K\}}$ in \mathcal{D}^K with $\tilde{\gamma} \ll \tilde{\gamma}'_a$ for all a ,

For all strategies ψ for all random variables Z taking values in $\{q\}$ and that are $\sigma(H_T)$ -measurable,

$$\sum_{a=1}^K E_\psi[N_a(T)] \text{KL}(\tilde{\gamma}_a, \tilde{\gamma}'_a) = \text{KL}(P_T^{H_T}, P_T^{H'_T}) \geq \text{KL}(\text{Ber}(E_\psi[Z]), \text{Ber}(E_{\psi'}[Z]))$$

The dependence on the strategy is hidden in the $E_\psi[N_a(T)]$, $E_\psi[Z]$ and $E_{\psi'}[Z]$

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 $K_{\text{inf}}(\vec{v}_a, \mu^*) = \inf \{ KL(\vec{v}_a, \vec{v}'_a) : \vec{v}'_a \in \mathcal{D} \text{ and } E(\vec{v}'_a) > \mu^* \}$

$$= \inf \{ KL(\vec{v}_a, \vec{v}'_a) : \vec{v}'_a \in \mathcal{D}, \vec{v}'_a \ll \vec{v}_a \text{ and } E(\vec{v}'_a) > \mu^* \}$$

(cf. convention: $\inf \emptyset = +\infty$ and the fact that $KL(\vec{v}_a, \vec{v}'_a) = +\infty$ when $\vec{v}'_a \ll \vec{v}_a$)

This is why we will

- Fix $\mathcal{D}, \Psi, \mathcal{T}$ and a st. $\Delta_a > 0$
- Fix an alternative model \vec{v}' of the form

$$\begin{cases} \vec{v}'_k = \vec{v}_k & \forall k \neq a \\ \vec{v}'_a & \text{s.t. } \vec{v}'_a \in \mathcal{D}, \vec{v}'_a \ll \vec{v}_a \text{ and } E(\vec{v}'_a) > \mu^* \end{cases}$$

That is, \vec{v} and \vec{v}' only differ at a ; a is the unique optimal arm in \vec{v}'

- Take $Z = N_a(\mathcal{T})/\mathcal{T}$ which is indeed $[\alpha_1]$ -valued $\sigma(H_T)$ -measurable

Our fundamental inequality yields, since \vec{v} and \vec{v}' only differ at a :

$$\begin{aligned} E_{\vec{v}}[N_a(\mathcal{T})] \cdot KL(\vec{v}_a, \vec{v}'_a) &\geq KL(Ber(E_{\vec{v}}[N_a(\mathcal{T})/\mathcal{T}]), Ber(E_{\vec{v}'}[N_a(\mathcal{T})/\mathcal{T}])) \\ &\geq -\ln 2 + (1 - E_{\vec{v}}[N_a(\mathcal{T})/\mathcal{T}]) \ln \frac{1}{1 - E_{\vec{v}'}[N_a(\mathcal{T})/\mathcal{T}]} \end{aligned}$$

Indeed: $KL(Ber(p), Ber(q))$

$$\begin{aligned} &= p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} \\ &= p \ln \frac{1}{q} + (1-p) \ln \frac{1}{1-q} + \underbrace{(p \ln p + (1-p) \ln(1-p))}_{\geq -\ln 2 \text{ by a simple function study over } [\alpha_1]} \\ &\geq -\ln 2 + (1-p) \ln \frac{1}{1-q} \end{aligned}$$

for all $p, q \in (\alpha_1)$ and even for all $p, q \in [\alpha_1]$ (study the cases $q=0$ and $q=1$ separately)

Now, the considered strategy Ψ is consistent and:

- in the problem \vec{v} , a is suboptimal: $E_{\vec{v}}[N_a(\mathcal{T})/\mathcal{T}] \rightarrow 0$

- in the problem \mathcal{S}^* , all arms $k \neq a$ are suboptimal:

$$\text{for all } \alpha \in (0,1], \quad T - E_{\mathcal{S}^*}[N_a(T)] = \sum_{k \neq a} E_{\mathcal{S}^*}[N_k(T)] = o(T^\alpha).$$

↳ in particular, for T large enough,

$$\frac{1}{1 - E_{\mathcal{S}^*}[N_a(T)/T]} = \frac{T}{T - E_{\mathcal{S}^*}[N_a(T)]} \geq \frac{T}{T^\alpha} = T^{1-\alpha}$$

Substituting back and dividing by $\ln T$:

for all $\alpha \in (0,1]$, for T large enough:

$$\frac{E_{\mathcal{S}^*}[N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}'_a)}{\ln T} \geq -\frac{\ln 2}{\ln T} + \left(1 - E_{\mathcal{S}^*}\left[\frac{N_a(T)}{T}\right]\right) \underbrace{\frac{\ln T^{1-\alpha}}{\ln T}}_{\rightarrow 0} = 1 - \alpha$$

thus

$$\liminf_{T \rightarrow +\infty} \frac{E_{\mathcal{S}^*}[N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}'_a)}{\ln T} \geq 1 - \alpha$$

Letting $\alpha \rightarrow 0$,

$$\liminf_{T \rightarrow +\infty} \frac{E_{\mathcal{S}^*}[N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}'_a)}{\ln T} \geq 1$$

Whether $\text{KL}(\bar{v}_a, \bar{v}'_a) < \infty$ or $= +\infty$, we thus get

$$\liminf_{T \rightarrow +\infty} \frac{E_{\mathcal{S}^*}[N_a(T)]}{\ln T} \geq \frac{1}{\text{KL}(\bar{v}_a, \bar{v}'_a)}$$

The left-hand side is independent of $\bar{v}'_a \in \mathcal{D}$ s.t. $\bar{v}'_a \gg \bar{v}_a$ and $E(\bar{v}'_a) \geq v^*$,

so that taking the supremum of the right-hand side over these \bar{v}'_a ,

we get the desired $1/\text{KL}(v_a, v^*)$ lower bound.

Proof of the lemma:

- The inequality \geq is a direct application of the data-processing inequality with expectations

and similarly $P_{\gamma}^{H_T} = K_T(K_{T-1}(\dots(K_1 m))\dots)$

- For the equality:

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

we check
below that
it is regular

{ where K_t is the transition kernel:

$$h \in [q] \times (\mathbb{R} \times [q])^{t-1} \mapsto K_t(h, \cdot) = \sum_{y_{t-1}(h)} \psi_{t-1}(y) \otimes m$$

Indeed: $T=0: H_0 = U_0 \sim \mathcal{U}_{[q]}$: $P_{\gamma}^{U_0} = m$

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

$$P_{\gamma}^{H_{t+1}}(A \times B' \times B) = P_{\gamma}(H_t \in A \text{ and } Y_{t+1} \in B' \text{ and } U_{t+1} \in B)$$

$$= \mathbb{E}_{\gamma} \left[\mathbb{1}_A(H_t) P_{\gamma}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B | H_t) \right]$$

tover
rule

definition
of the bandit
model and
of the strategy

$$\mathbb{E}_{\gamma} \left[\mathbb{1}_A(H_t) \sum_{y_t(H_t)} \psi_t(y_t) (B') m(B) \right]$$

) definition of
 K_{t+1}

$$= \mathbb{E}_{\gamma} \left[\mathbb{1}_A(H_t) K_{t+1}(H_t, B' \times B) \right]$$

) were rewriting

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

) definition of
 $K_{t+1} P_{\gamma}^{H_t}$

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

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

- The K_t are regular transition kernels: $\forall E \in \mathcal{B}(\mathbb{R}) \cap \mathcal{B}([q])$,

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

is measurable as ψ_t is measurable

- Assumption (*): $\forall h, K_t(h, \cdot) \ll K_t'(h, \cdot)$ as $K_t \ll \tilde{v}_a \ll v_a$ by assumption

* Assumption (**):

$$(h, (y, u)) \mapsto \frac{dK_t(h, \cdot)}{dK_t^*(h, \cdot)} (y, u)$$

is indeed bi-measurable
(product of measurable functions)

$$= \sum_{\alpha=1}^K \prod_{\gamma=1}^{\alpha} \frac{d\zeta_\alpha}{d\psi_t(h) = \alpha} \frac{d\zeta_\alpha}{d\psi_t^*(h)} (y)$$

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

- $KL(P_{y^0}^{H^0}, P_{y^0}^{H^0}) = KL(\eta, \eta) = 0$

- For $t > 0$,

$$\begin{aligned} & KL(P_{y^t}^{H^t}, P_{y^t}^{H^t}) \\ &= KL(K_{t+1} P_{y^t}^{H^t}, K_{t+1}' P_{y^t}^{H^t}) \\ &= KL(P_{y^t}^{H^t}, P_{y^t}^{H^t}) + \int KL(K_{t+1}(h, \cdot), K_{t+1}'(h, \cdot)) dP_{y^t}^{H^t}(h) \\ &= KL(P_{y^t}^{H^t}, P_{y^t}^{H^t}) + \int KL(\underbrace{\zeta_t(h) \otimes \eta}_{\cancel{\zeta_t(h) \otimes \eta}}, \underbrace{\zeta_t'(h) \otimes \eta}_{\cancel{\zeta_t'(h) \otimes \eta}}) dP_{y^t}^{H^t}(h) \\ &= KL(P_{y^t}^{H^t}, P_{y^t}^{H^t}) + \sum_{\alpha=1}^K \underbrace{KL(\zeta_\alpha, \zeta_\alpha')}_{\int \mathbb{1}_{\{\psi_t(h)=\alpha\}} dP_{y^t}^{H^t}(h)} \underbrace{\mathbb{E}\left[\mathbb{1}_{\{\psi_t(h)=\alpha\}}\right]}_{\mathbb{E}\left[\mathbb{1}_{\{\psi_t(H^t)=\alpha\}}\right]} \end{aligned}$$

Exercise:

$$\frac{1}{K \text{inf}(\bar{x}_i, \mu^*) \Delta_i} \quad \text{vs.} \quad \frac{8}{\Delta_i^2} \quad \text{for UCB}$$

Recall that in the model $\mathcal{D} = \mathcal{P}(\mathcal{A}|I)$, the UCB algorithm employs the following performance bound:

$$\forall i \in \mathcal{P}(\mathcal{A}|I)^K, \quad \forall \alpha \text{ s.t. } \Delta_i > 0,$$

$$E_{\pi} [N_i(T)] \leq \frac{8}{\Delta_i^2} \ln T + 2.$$

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

But how do these $\frac{8}{\Delta_i^2}$ and $\frac{2}{\Delta_i^2}$ constants compare to $\frac{1}{K \text{inf}(\bar{x}_i, \mu^*) \Delta_i}$?

(1) For $p, q \in [0, 1]$, we denote

$$kl(p|q) = KL(Ber(p), Ber(q))$$

Show that $\forall (p|q) \in [0, 1]^2, \quad kl(p|q) \geq 2(p-q)^2$.

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

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

The total variation distance between P and Q

Even better, show the stronger form: $\sup_{Z \in \mathcal{F}\text{-measurable taking values in } [0, 1]} |E_P[Z] - E_Q[Z]| \leq \sqrt{\frac{1}{2} KL(P, Q)}$

(3) Exhibit a lower bound on $K \text{inf}(\bar{x}_i, \mu^*) \Delta_i$ 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:

- $\forall 0 \leq p < q \leq 1$,

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

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

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

(2) Show that all strategies smarts than the uniform strategy [ie, such that for all bandit problems $\forall a$ s.t. $\mu_a = \mu^*$, $E[N_a(T)] \geq T_K$], we have:

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

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

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

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

at
least half
of the time

uniform exploration

Hint: Consider the same alternative bandit problems as in the theorem giving the asymptotic lower bound.