

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 follow 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 $\forall B \in \mathcal{F}'$ $\omega \mapsto K(\omega, B)$ 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?)

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

$$KP(A \times B) = \int_{\Omega} \underbrace{\mathbf{1}_A(w)}_{\text{is } \mathcal{F}\text{-measurable}} K(w, B) dP(w)$$

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

An extension of
Fubini (-Tonelli) theorem
↓

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

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

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

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 KP

(sketch) Extension to $\mathbf{1}_E$ for any $E \in \mathcal{F} \times \mathcal{F}'$ by an argument of σ -algebra contained / monotone class theorem, using (monotone) convergence (including the $w \mapsto \int \varphi(w, \cdot) K(w, d\cdot)$ measurability)

Extension to $\varphi \geq 0$ by monotone convergence

Extension to $\varphi \in L^1$ by dominated convergence

Question: Does anyone have a simpler argument?

Theorem [chain rule for KL]:

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

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

Then

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

where

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

Remark: see last page of these lecture notes 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

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

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

We will not if we use this, actually

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

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

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

By Radon-Nikodym's Theorem:

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

and the converse is easily seen

* Therefore, we have $KIP \ll LQ$ as soon as $P \ll Q$, $KIP \ll LQ$ and $P \ll Q$ are thus assumed with no loss of generality that putative equality is $+\infty = +\infty$.

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

$\Psi = \int g \ln(g)$ is lower bounded, Pfeffer lemma (extension of Fubini-Tonelli) extends to it:

$$\begin{aligned} & \int g \ln(g) d\mathbb{Q} \\ &= \int_{\Omega} \left(\left(\int_{\Omega'} g(w, w') \ln g(w, w') L(w, dw') \right) + \ln f(w) \right) f(w) d\mathbb{Q}(w) \\ &= \int_{\Omega} \left(f(w) \underbrace{\text{KL}(K(w), L(w))}_{\geq 0} + \underbrace{f(w) \ln f(w)}_{\text{lower bounded}} \right) d\mathbb{Q}(w) \end{aligned}$$

$\sum_{i=1}^n$ sum of
two bounded
functions

$$\int_{\Omega} \text{KL}(K(w), L(w)) \underbrace{f(w) d\mathbb{Q}(w)}_{dP(w)} + \underbrace{\int f \ln f d\mathbb{Q}}_{\text{KL}(P, Q)}$$

REMARKS ON THE ASSUMPTIONS.

- The assumptions (i) and (ii) 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 $(O_n)_{n \in \mathbb{N}}$ 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 $(O_n)_{n \in \mathbb{N}}$.

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: Marin Brilu + Hedi Thalibi, M2 students of Spring 2017

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 $\Delta_a > 0$

$U_0, U_1, U_2, \dots, U_{[T]}$

decision-maker
but does not
know specific
arm i

- 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 $\pi_{I_t}^*$ (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 $\mathcal{T} \in \mathbb{S}^K$,

$$\forall \delta \in (0, 1], \quad \forall \alpha \text{ s.t. } \Delta_\alpha > 0, \quad \mathbb{E}[N_a(T)] = o(T^\alpha).$$

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

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

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

$$\forall \mathcal{T} \in \mathbb{S}^K, \quad \forall \alpha \text{ s.t. } \Delta_\alpha > 0,$$

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

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

- Optimal (in some sense) such constant: $C_{\text{opt}}(\mathcal{T}) = \frac{1}{K_{\text{inf}}(\bar{v}_a, \mu^*, \mathcal{D})} = \frac{1}{K_{\text{inf}}(\bar{v}_a, \mu^*)}$

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

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

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

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 $\mathcal{T} = (\bar{v}_a)_{a \in [1, \dots, K]} \in \mathbb{S}^K$,

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

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

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

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\ln T} \geq \sum_{a: A_a > 0} \frac{\Delta_a}{\text{KL}_{\psi}(\vec{\nu}_a, \mu^*_a, \mathcal{D})}.$$

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

Example: $P_\psi^{H_T}$ 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, $P_\psi^{H_T}$ 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 $\vec{\nu} = (\vec{\nu}_a)_{a \in [1..K]}$ and $\vec{\nu}' = (\vec{\nu}'_a)_{a \in [1..K]}$ in \mathcal{D}^K with $\vec{\nu}_a \ll \vec{\nu}'_a$ for all a ,

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

$$\begin{aligned} \sum_{a=1}^K E_{\vec{\nu}}[\ln_a(T)] \text{KL}(\vec{\nu}_a, \vec{\nu}'_a) &= \text{KL}(P_{\vec{\nu}}^{H_T}, P_{\vec{\nu}'}^{H_T}) \\ &\geq \text{KL}(\text{Ber}(E_{\vec{\nu}}[Z]), \text{Ber}(E_{\vec{\nu}'}[Z])) \end{aligned}$$

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 $K_{\text{inf}}(\tilde{\nu}_a, \mu^*) = \inf \{ \text{KL}(\tilde{\nu}_a, \tilde{\nu}_a^i) : \tilde{\nu}_a^i \in \mathcal{D} \text{ and } E(\tilde{\nu}_a^i) > \mu^* \}$

$$= \inf \{ \text{KL}(\tilde{\nu}_a, \tilde{\nu}_a^i) : \tilde{\nu}_a^i \in \mathcal{D}, \tilde{\nu}_a^i \ll \tilde{\nu}_a^i \text{ and } E(\tilde{\nu}_a^i) > \mu^* \}$$

(cf. convention: $\inf \emptyset = +\infty$ and the fact that $\text{KL}(\tilde{\nu}_a, \tilde{\nu}_a^i) = +\infty$ when $\tilde{\nu}_a \ll \tilde{\nu}_a^i$)

This is why we will

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

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

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

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

Our fundamental inequality yields, since $\tilde{\nu}$ and $\tilde{\nu}^*$ only differ at a :

$$\begin{aligned} E_{\tilde{\nu}}[N_a(\mathcal{T})] \text{KL}(\tilde{\nu}_a, \tilde{\nu}_a^*) &\geq \text{KL}\left(\text{Ber}\left(E_{\tilde{\nu}}[N_a(\mathcal{T})/\mathcal{T}]\right), \text{Ber}\left(E_{\tilde{\nu}^*}[N_a(\mathcal{T})/\mathcal{T}]\right)\right) \\ &\geq -\ln 2 + \left(1 - E_{\tilde{\nu}}[N_a(\mathcal{T})/\mathcal{T}]\right) \ln \frac{1}{1 - E_{\tilde{\nu}^*}[N_a(\mathcal{T})/\mathcal{T}]} \end{aligned}$$

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

$$\begin{aligned} &= p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} \\ &= \underbrace{p \ln \frac{1}{q}}_{\geq 0} + (1-p) \ln \frac{1}{1-q} + \underbrace{(p \ln p + (1-p) \ln (1-p))}_{\geq -\ln 2 \text{ by a simple function study over } [\alpha_1]} \\ &\geq -\ln 2 + (1-p) \ln \frac{1}{1-q} \end{aligned}$$

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 $\tilde{\nu}$, a is suboptimal: $E_{\tilde{\nu}}[N_a(\mathcal{T})/\mathcal{T}] \rightarrow 0$

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

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

↳ in particular, for T large enough,

$$\frac{1}{1 - \mathbb{E}_{\mathcal{V}^*} [N_a(T)/T]} = \frac{T}{T - \mathbb{E}_{\mathcal{V}^*} [N_a(T)]} \geq \frac{T}{T^\alpha} = T^{1-\alpha}$$

Substituting back and dividing by $\ln T$:

$$\frac{\mathbb{E}_{\mathcal{V}^*} [N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}_a^*)}{\ln T} \geq -\frac{\ln 2}{\ln T} + \left(1 - \mathbb{E}_{\mathcal{V}^*} \left[\frac{N_a(T)}{T} \right] \right) \underbrace{\frac{\ln T^{1-\alpha}}{\ln T}}_{\rightarrow 0} = 1 - \alpha$$

thus

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{V}^*} [N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}_a^*)}{\ln T} \geq 1 - \alpha$$

$$\text{Letting } \alpha \rightarrow 0, \quad \liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{V}^*} [N_a(T)]}{\ln T} \frac{\text{KL}(\bar{v}_a, \bar{v}_a^*)}{\ln T} \geq 1$$

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

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_{\mathcal{V}^*} [N_a(T)]}{\ln T} \geq \frac{1}{\text{KL}(\bar{v}_a, \bar{v}_a^*)}$$

The left-hand side is independent of $\bar{v}_a^* \in \mathcal{S}$ s.t. $\bar{v}_a^* \gg \bar{v}_a$ and $E(\bar{v}_a^*) \geq \mu^*$,

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

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

Proof of the lemma: • The inequality \geq is a direct application of the data-processing inequality with expectations, see the previous lecture for its statement.

- For the equality: We will explain how $\mathbb{P}_{\gamma}^{H_T}$ is constructed.

With no loss of generality, we can consider that

- the underlying probability space is $\Omega = [0,1] \times (\mathbb{R} \times [0,1])^T$
- H_T is the identity over Ω , i.e., that the $U_0, Y_1, U_1, \dots, Y_T, U_T$ are the projections on each component,
- \mathbb{P}_{γ} is given by

$$\forall B \in \mathcal{B}([0,1]), \quad \mathbb{P}_{\gamma}(U_0 \in B) = \eta(B)$$

$$\forall t \in \{0, \dots, T-1\}, \quad \forall B' \in \mathcal{B}(\mathbb{R}), \quad \forall B \in \mathcal{B}([0,1]), \quad \mathbb{P}_{\gamma}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B \mid H_t) = \int_{\psi_t(H_t)} (B') \eta(B)$$

where $\mathcal{B}(S)$ is the Borel- σ -algebra of a set $S \subseteq \mathbb{R}$
 η is the Lebesgue measure over $[0,1]$

In particular: $\mathbb{P}_{\gamma}^{H_t}$ refers to the first $T-t$ marginals of $\mathbb{P}_{\gamma}^{H_T} = \mathbb{P}_{\gamma}^{H_t}$.

A similar construction can be done for the bandit problem γ' .

Now, the equality $\mathbb{P}_{\gamma}^{H_t}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B \mid H_t) = \int_{\psi_t(H_t)} (B') \eta(B)$

indicates that $\mathbb{P}_{\gamma}^{H_{t+1}} = K_t \mathbb{P}_{\gamma}^{H_t}$ for the regular transition kernel

details on next page ↑ regularly is:
for $E \in \mathcal{B}(\mathbb{R}) \otimes \mathcal{B}([0,1])$
 $h \mapsto \psi_t(h, E)$ is
measurable,
it follows from the
measurability of ψ_t

Let us check the assumptions of our chain rule:

$$(*) \quad \forall h_i, \quad K_t(h_i) \ll K'_t(h_i) \quad \text{as } \forall a, \quad \nu_a \ll \nu'_a \text{ by assumption}$$

$$(**) \quad (h, (y_u)) \mapsto \frac{dK_t(h_i)}{dK'_t(h_i)}(y_u) = \sum_{a=1}^K \mathbb{1}_{\{\psi_t(h)=a\}} \frac{d\nu_a}{d\nu'_a}(y) \\ \text{is indeed bi-measurable.}$$

Therefore, for $t \in \{0, \dots, T-1\}$,

$$\begin{aligned} \text{KL}(\mathbb{P}_{\gamma}^{H_{t+1}}, \mathbb{P}_{\gamma'}^{H_{t+1}}) &= \text{KL}(\mathbb{P}_{\gamma}^{H_t}, \mathbb{P}_{\gamma'}^{H_t}) + \int \text{KL}(\nu_{\psi_t(h)} \ll \nu'_{\psi_t(h)}, \nu'_{\psi_t(h)} \ll \nu_{\psi_t(h)}) \, d\mathbb{P}_{\gamma}^{H_t}(h) \\ &= \text{KL}(\mathbb{P}_{\gamma}^{H_t}, \mathbb{P}_{\gamma'}^{H_t}) + \sum_{a=1}^K \text{KL}(\nu_a, \nu'_a) \, \mathbb{P}_{\gamma}^{H_t} \{ \psi_t(h) = a \} \end{aligned}$$

Now, $I_{t+1} = \Psi_t(H_t)$, so that

$$\begin{aligned} P_{\gamma}^{H_t} \{ \Psi_t(h) = a \} &= P_{\gamma}^{\Psi_t^0} \{ \Psi_t(h) = a \} \\ &= P_{\gamma} \{ I_{t+1} = a \} = \mathbb{E}_{\gamma} \left[\frac{1}{I_{t+1}} \right] \end{aligned}$$

Summing up:

$$- \text{KL}(P_{\gamma}^{H_0}, P_{\gamma}^{H_0}) = \text{KL}(P_{\gamma}^{U_0}, P_{\gamma}^{U_0}) = \text{KL}(\eta, \eta) = 0$$

$$\begin{aligned} - \forall t \in \{0, \dots, T-1\}, \quad \text{KL}(P_{\gamma}^{H_{t+1}}, P_{\gamma}^{H_{t+1}}) &= \text{KL}(P_{\gamma}^{H_t}, P_{\gamma}^{H_t}) \\ &\quad + \sum_{a=1}^K \text{KL}(\delta_a, \delta_a) \mathbb{E}_{\gamma} \left[\frac{1}{I_{t+1}} \right] \end{aligned}$$

so that the stated result follows by induction.

*Explanation of why $P_{\gamma}^{H_t} = K_t P_{\gamma}^{H_t}$:

$$\forall A \in \mathcal{B}([q]) \times ((R \times [q]))^t, \quad \forall B \in \mathcal{B}(R), \quad B \subseteq \mathcal{B}([q]),$$

$$K_t P_{\gamma}^{H_t} (A \times (B \setminus B))$$

$$\stackrel{\text{def of product}}{=} \int_{[q] \times ((R \times [q]))^t} \mathbb{1}_A(h) K_t(h, B \setminus B) dP_{\gamma}^{H_t}(h)$$

$$\stackrel{\text{def of image distribution}}{=} \mathbb{E}_{\gamma} \left[\mathbb{1}_A(H_t) K_t(H_t, B \setminus B) \right]$$

$$\stackrel{\text{K}_t \text{ defined as joint conditional probability}}{=} \mathbb{E}_{\gamma} \left[\mathbb{1}_A(H_t) P_{\gamma} \left(Y_{t+1} \in B' \text{ and } U_{t+1} \in B \mid H_t \right) \right]$$

$$\stackrel{\text{tower rule}}{=} P_{\gamma} \left(H_t \in A \text{ and } Y_{t+1} \in B' \text{ and } U_{t+1} \in B \right) = P_{\gamma}^{H_t} (A \times B' \times B).$$

Exercise:

$$\frac{1}{K \text{inf}(\mu_1, \mu^*, \Delta)} \quad \text{vs.} \quad \frac{8}{\Delta^2} \quad \text{for UCB}$$

Recall that in the model $\mathcal{D} = \mathcal{P}(\mathcal{G}_1)$, the UCB algorithm enjoys the following performance bound:

$$\forall i \in \mathcal{P}(\mathcal{G}_1)^K, \quad \forall t \text{ s.t. } \Delta > 0,$$

$$E_{\mathcal{P}}[N_{\mu}(t)] \leq \frac{8}{\Delta^2} \ln T + 2.$$

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

But how do these $\frac{8}{\Delta^2}$ and $\frac{2}{\Delta^2}$ constants compare to $\frac{1}{K \text{inf}(\mu_1, \mu^*, \mathcal{P}(\mathcal{G}_1))}$?

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

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

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

(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: $\sup_{Z \text{ F-measurable
taking values
in } [0, 1]} |\mathbb{E}_P[Z] - \mathbb{E}_Q[Z]| \leq \sqrt{\frac{1}{2} KL(P, Q)}$

(3) Exhibit a lower bound on $K \text{inf}(\mu_1, \mu^*, \mathcal{P}(\mathcal{G}_1))$ 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 Finsler's inequality:

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

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

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

(2) Show that all strategies smoothen 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 $K^{**} = \max_{j: \alpha_j > 0} K_{\inf}(j, \mu_j^*)$

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

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

at least half
of the time

uniform exploration

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

Distribution-free (ie uniform) lower bounds.

We prove: For all $K \geq 2$ and $T \geq K/5$

$$\inf_{\text{strategies } \psi} \sup_{\substack{\gamma_1, \dots, \gamma_K \\ \text{in } \mathcal{P}(\{1\})}} R_T \geq \frac{1}{20} \sqrt{T K}.$$

and even:

$\sup_{\substack{\gamma_1, \dots, \gamma_K \\ \text{being Bernoulli distributions}}} \dots$

(In class I did not write the correct $\gamma^{(i)}$ and $\gamma^{(j)}$)
it's with $\varepsilon_{1/2}$ in γ)

Proof: • We consider the bandit problem $\gamma^{(i)} = (\text{Ber}(\frac{1}{2}, \varepsilon_{1/2}), \dots, \text{Ber}(\frac{1}{2}, \varepsilon_{1/2}))$

versus the bandit problems $\gamma^{(i)} = (\text{Ber}(\frac{1}{2}, \varepsilon_{1/2}), \dots, \text{Ber}(\frac{1}{2}, \varepsilon_{1/2}), \text{Ber}(\frac{1}{2} + \varepsilon_i), \text{Ber}(\frac{1}{2} - \varepsilon_i), \dots, \text{Ber}(\frac{1}{2}, \varepsilon_{1/2}))$
for $i \in \{1..K\}$ and $\varepsilon \in (0, \frac{1}{2})$

For all strategies,

$$\sup_{\gamma \in \mathcal{P}(\{1\})} R_T \geq \sup_{\varepsilon \in (0, \frac{1}{2})} \max_{i \in \{1..K\}} \bar{R}_i$$

$$= \sup_{\varepsilon \in (0, \frac{1}{2})} \max_{i \in \{1..K\}} \sum_{k \neq i} \varepsilon \mathbb{E}_{\gamma^{(i)}} [N_k(T)]$$

↑ expected number
of times k
is pulled under $\gamma^{(i)}$

$$\text{And } \sum_{k \neq i} N_k(T) = T - N_i(T)$$

$$\text{Thus, } \sup_{\gamma \in \mathcal{P}(\{1\})} R_T \geq \sup_{\varepsilon \in (0, \frac{1}{2})} \max_{i \in \{1..K\}} \mathbb{E} (T - \mathbb{E}_{\gamma^{(i)}} [N_i(T)])$$

• Now, we use the fundamental inequality to upper bound one of the $\mathbb{E}_{\gamma^{(i)}} [N_i(T)]$

There exists $k \in \{1..K\}$ such that $\mathbb{E}_{\gamma^{(k)}} [N_k(T)] \leq T/K$. For this k :

$$\sum_{a=1}^K \mathbb{E}_{\gamma^{(k)}} [N_a(T)] \text{KL}(\gamma^{(k)}, \gamma^{(k)}) \geq k \mathbb{E}_{\gamma^{(k)}} [N_k(T)/T] \mathbb{E}_{\gamma^{(k)}} [N_k(T)/T]$$

Since $\gamma^{(k)}$ and $\gamma^{(k)}$ only differ at arm k ,

$$= \mathbb{E}_{\gamma^{(k)}} [N_k(T)] \text{KL}(\text{Ber}(\frac{1}{2}, \varepsilon_{1/2}), \text{Ber}(\frac{1}{2} + \varepsilon_{1/2}))$$

$$\leq \frac{1}{K} k \text{KL}(\frac{1}{2} - \varepsilon_{1/2}, \frac{1}{2} + \varepsilon_{1/2})$$

by definition
of k

$\text{KL} = \text{KL}$ between
Bernoulli distributions

we lower bound this side
by Pinsker's inequality for
Bernoulli distributions:

$$\geq 2 \left(\mathbb{E}_{\gamma^{(k)}} \left[\frac{N_k(T)}{T} \right] - \mathbb{E}_{\gamma^{(k)}} \left[\frac{N_k(T)}{T} \right]^2 \right)^2$$

We proved so far:

$$2 \left(\mathbb{E}_{\gamma^{(k)}} \left[\frac{N_k(T)}{T} \right] - \mathbb{E}_{\gamma^{(k)}} \left[\frac{N_k(T)}{T} \right]^2 \right)^2 \leq \frac{T}{K} k \left(\frac{1}{2} - \varepsilon_{1/2}, \frac{1}{2} + \varepsilon_{1/2} \right)$$

This entails in particular :

$$\mathbb{E}_{\underline{\omega}^{(t)}} \left[\frac{N_k(t)}{T} \right] \leq \mathbb{E}_{\underline{\omega}^{(t)}} \left[\frac{N_k(t)}{T} \right] + \underbrace{\sqrt{\frac{T}{2K} k l(\frac{1}{2} - \varepsilon_k, \frac{1}{2} + \varepsilon_k)}}_{\leq \gamma_k \text{ by definition of } k}$$

- Substituting in the regret lower bound we get (with $i=k$) :

$$\sup_{\underline{\omega} \in J([q_1])} \bar{R}_T \geq \sup_{\mathcal{E} \in (\frac{1}{2}, 1)} \mathcal{E} T \left(1 - \mathbb{E}_{\underline{\omega}^{(t)}} \left[\frac{N_k(t)}{T} \right] \right)$$

$$\geq \sup_{\mathcal{E} \in (\frac{1}{2}, 1)} \mathcal{E} T \left(1 - \frac{\gamma_k}{\sqrt{\frac{T}{2K}}} - \sqrt{\frac{T}{2K} k l(\frac{1}{2} - \varepsilon_k, \frac{1}{2} + \varepsilon_k)} \right)$$

$$\geq \sup_{\mathcal{E} \in (\frac{1}{2}, 1)} \mathcal{E} T \left(\frac{1}{2} - \sqrt{\frac{T}{2K} k l(\frac{1}{2} - \varepsilon_k, \frac{1}{2} + \varepsilon_k)} \right)$$

$$\text{Now, } k l(\frac{1}{2} - \varepsilon_k, \frac{1}{2} + \varepsilon_k) = (\frac{1}{2} - \varepsilon_k) \ln \frac{1-\varepsilon}{1+\varepsilon} + (\frac{1}{2} + \varepsilon_k) \ln \frac{1+\varepsilon}{1-\varepsilon} \\ = \varepsilon \ln \frac{1+\varepsilon}{1-\varepsilon}$$

$$\leq 2.5 \varepsilon^2 \text{ on } (0, \frac{1}{2}) \\ (\text{study } \varepsilon \mapsto \varepsilon \ln(\frac{1+\varepsilon}{1-\varepsilon})/\varepsilon^2)$$

So that the regret lower bound becomes

$$\sup_{\underline{\omega} \in J([q_1])} \bar{R}_T \geq \sup_{\mathcal{E} \in (\frac{1}{2}, 1)} \frac{T}{2} \left(\mathcal{E} - \mathcal{E} \sqrt{\frac{5T}{K}} \right)$$

$$\text{we pick } \mathcal{E} \text{ st. } 1 - 2\mathcal{E} \sqrt{\frac{5T}{K}} = 0$$

$$\text{that is, } \mathcal{E} = \sqrt{\frac{K}{5}} / 2\sqrt{5}$$

which is indeed $\frac{1}{2}$ as soon as $T > K/5$.

We get a $\sqrt{K} \left(\frac{1}{2} \left(\frac{1}{2}\sqrt{5} - \frac{1}{20}\sqrt{5} \right) \right) \geq \frac{1}{20} \sqrt{K}$ bound.

Note: The original proof used $\text{Ber}(\frac{1}{2})$ and $\text{Ber}(\frac{1}{2} + \varepsilon)$ but the approach with $\text{Ber}(\frac{1}{2} - \varepsilon)$ and $\text{Ber}(\frac{1}{2} + \varepsilon)$ leads to slightly simpler calculation.