Let's first complete the proof of the Lemma:      [ "Hoeffding-Azuma inequality with a random number of summands" ]

Setting: Probability distributions $\nu_1, \dots \nu_K$ over $[0,1]$

with respective expectations $\mu_1, \dots \mu_K$

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

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

ie: $y_t \mid I_t \sim \nu_{I_t}$

We denote $N_a(t) = \sum\limits_{s=1}^{t} \mathbb{1}_{\{I_s = a\}}$ and assume that each arm $a$ was pulled once in the first $K$ rounds,

so that: $N_a(t) \geqslant 1 \quad \forall t \geqslant K$

Then, for $t \geqslant K$: $\hat{\mu}_{a,t} = \dfrac{1}{N_a(t)} \sum\limits_{s=1}^{t} y_s \, \mathbb{1}_{\{I_s = a\}}$

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

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

$$z_t = \sum\limits_{s=1}^{t} (y_s - \mu_a) \, \mathbb{1}_{\{I_s = a\}}$$

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

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

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

where we used the bandit model $= \left( \mu_{I_t} - \mu_a \right) \mathbb{1}_{\{I_t = a\}} = 0 \quad \text{a.s.}$

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

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

$a_t = -\mu_a$ and $b_t = 1 - \mu_a$

so that

(actually in the version I stated, I can have $\leqslant$ or $<$) $(b_t - a_t)^2 = 1$

$$1 - t\delta \leq \mathbb{P}\left\{ Z_t < \sqrt{\frac{t}{2}\ln\frac{1}{t\delta}} \right\} = \mathbb{P}\left\{ N_a(t)\left(\hat{\mu}_{a,t} - \mu_a\right) < \sqrt{\frac{t}{2}\ln\frac{1}{t\delta}} \right\}$$

$$= \mathbb{P}\left\{ \hat{\mu}_{a,t} - \sqrt{\frac{t}{N_1(t)}}\sqrt{\frac{\ln(1/t\delta)}{2\,N_a(t)}} < \mu_a \right\}$$

versus the bound of our lemma:  $\quad 1 - t\delta \leq \mathbb{P}\left\{ \hat{\mu}_{a,t} - \sqrt{\frac{\ln(1/\delta)}{2 N_1(t)}} < \mu_a \right\}$

The proposed deviations essentially differ from a $\sqrt{t/N_1(t)}$ factor, and it is so nice to get rid of it!

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

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

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

but we can do better!

$$\mathbb{E}\left[ e^{x(Y_t - \mu_a)\,\mathbb{1}_{\{I_t = a\}}} \mid \mathcal{F}_{t-1} \right] \leq e^{x^2/8} \quad a.s.$$

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

$$\mathbb{E}\left[ e^{x(Y_t - \mu_a)\,\mathbb{1}_{\{I_t = a\}}} \mid \mathcal{F}_{t-1} \right] = \mathbb{E}\left[ e^{x(Y_t - \mu_a)\,\mathbb{1}_{\{I_t = a\}}}\left( \mathbb{1}_{\{I_t = a\}} + \mathbb{1}_{\{I_t \neq a\}} \right) \mid \mathcal{F}_{t-1} \right]$$

given what we had before

$$= \mathbb{E}\left[ e^{x(Y_t - \mu_a)\,\mathbb{1}_{\{I_t = a\}}} \mid \mathcal{F}_{t-1} \right] \mathbb{1}_{\{I_t = a\}} + e^0\,\mathbb{1}_{\{I_t \neq a\}}$$

$$\leq e^{x^2/8}\,\mathbb{1}_{\{I_t = a\}} + \mathbb{1}_{\{I_t \neq a\}} = \exp\left( \frac{x^2}{8}\,\mathbb{1}_{\{I_t = a\}} \right)$$

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

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

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

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

(2) We prove that $\forall \varepsilon > 0, \quad \forall \ell \geq 1,$

$$\mathbb{P}\{Z_t \geq \varepsilon \text{ and } N_a(t) = \ell\}$$
$$\leq \exp(-2\varepsilon^2/\ell)$$

Indeed, by a Markov-Chernov bounding,

$\forall x > 0,$

$$\mathbb{P}\{Z_t \geq \varepsilon \text{ and } N_a(t) = \ell\} \leq e^{-x\varepsilon} \mathbb{E}\left[e^{xZ_t} \mathbb{1}_{\{N_a(t) = \ell\}}\right]$$

$$= e^{-x\varepsilon + \frac{x^2\ell}{8}} \mathbb{E}\left[e^{xZ_t - \frac{x^2}{8}N_a(t)} \mathbb{1}_{\{N_a(t) = \ell\}}\right]$$

$$\leq e^{-x\varepsilon + x^2\ell/8} \underbrace{\mathbb{E}\left[e^{xZ_t - \frac{x^2}{8}N_a(t)}\right]}_{\leq 1 \text{ by } (1)}$$

Optimizing over $x > 0$

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

(3) Conclusion: we prove that $\mathbb{P}\left\{\mu_a \leq \hat{\mu}_{a,t} - \sqrt{\dfrac{\ln(1/\delta)}{2N_a(t)}}\right\} \leq t\delta$

Indeed, by the union bound,

$$\mathbb{P}\left\{\mu_a \leq \hat{\mu}_{a,t} - \sqrt{\dfrac{\ln(1/\delta)}{2N_a(t)}}\right\}$$

$$= \sum_{\ell=1}^{t} \mathbb{P}\left\{N_a(t) = \ell \text{ and } \mu_a \leq \hat{\mu}_{a,t} - \sqrt{\ln(1/\delta)/2\ell}\right\}$$

$$= \sum_{\ell=1}^{t} \mathbb{P}\left\{N_a(t) = \ell \text{ and } \dfrac{Z_t}{N_a(t)} \geq \sqrt{\ln(1/\delta)/2\ell}\right\}$$

$$= \sum_{\ell=1}^{t} \mathbb{P}\left\{N_a(t) = \ell \text{ and } Z_t \geq \sqrt{\ell \ln(1/\delta)/2}\right\}$$

by (2) $\Bigg( \qquad \leq \sum_{\ell=1}^{t} \exp\left(-2\left(\ell \ln(1/\delta)/2\right)/\ell\right) = t\delta.$

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

Assume that $\nu_1, \dots \nu_K$ are unknown but that the decision-maker knows $\nu_j \in \mathscr{D}$ $\forall j$.

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

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

$$\overline{R}_T \text{ is at best of the order of } \left(\sum_{a:\Delta_a > 0} \frac{\Delta_a}{K_{\inf}(\nu_{a}, \mu^*, \mathscr{D})}\right) \ln T$$

where

$$K_{\inf}(\nu_a, \mu^*, \mathscr{D}) = \inf\left\{ KL(\nu_a, \nu_a') : \begin{array}{l} \nu_a' \in \mathscr{D} \\ \mathbb{E}(\nu_a') > \mu^* \end{array} \right\}$$

Kullback–        expectation
Leibler          of $\nu_a'$
divergence

We will do so by
 - proving a universal lower bound

 - exhibiting a strategy, called KL-UCB), to achieve the bound.

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

The    Kullback–Leibler divergence:    definition    and    basic properties.

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

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

Basic facts:

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

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

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

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

with equality if and only
if $\dfrac{d\mathbb{P}}{d\mathbb{Q}}$ is $\mathbb{Q}$-as constant, ie, $\mathbb{P} = \mathbb{Q}$

Exercise 1:    A useful rewriting.    Prove the following result:
Assume $\mathbb{P} \ll \mathbb{Q}$ and let $\gamma$ be any probability measure over $(\Omega, \mathcal{F})$
such that $\mathbb{P} \ll \gamma$ and $\mathbb{Q} \ll \gamma$. Denote $f = \dfrac{d\mathbb{P}}{d\gamma}$ and $g = \dfrac{d\mathbb{Q}}{d\gamma}$.

Then:    $KL(\mathbb{P}, \mathbb{Q}) = \int_{\Omega} \dfrac{f}{g} \ln\left( \dfrac{f}{g} \right) g \, d\mu$

$$= \int_{\Omega} \ln\left( \dfrac{f}{g} \right) f \, d\mu$$

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

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

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

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

Denote by $\mathbb{P}^X$ and $\mathbb{Q}^X$ the laws of $X$ under $\mathbb{P}$ and $\mathbb{Q}$.

Then :

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

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

$$\mathbb{P}^X \ll \mathbb{Q}^X, \quad \text{with} \quad \frac{d\mathbb{P}^X}{d\mathbb{Q}^X} = \mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X = \cdot\right] \overset{\text{not.}}{=} \gamma$$

ie, $\gamma(x) = \mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X\right]$.

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

$$\mathbb{P}^X(B) = \mathbb{P}\{X \in B\} = \int_{\Omega} \mathbb{1}_B(X) \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{Q} = \int_{\Omega} \mathbb{1}_B(X) \mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X\right] d\mathbb{Q}$$

$$\overset{\text{not.}}{=} \int_{\Omega} \mathbb{1}_B(X) \gamma(X) d\mathbb{Q} = \int_{\Omega'} \mathbb{1}_B \gamma \, d\mathbb{Q}^X.$$

by definition of $\mathbb{Q}^X$

Therefore,

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

definition of $\gamma$

$$= \int_{\Omega} \left(\mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X\right] \ln \mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X\right]\right) d\mathbb{Q}$$

conditional version of Jensen's inequality

$$\leq \int_{\Omega} \mathbb{E}_{\mathbb{Q}}\left[\frac{d\mathbb{P}}{d\mathbb{Q}} \ln \frac{d\mathbb{P}}{d\mathbb{Q}} \Big| X\right] d\mathbb{Q}$$

tower rule

$$= \int_{\Omega} \left(\frac{d\mathbb{P}}{d\mathbb{Q}} \ln \frac{d\mathbb{P}}{d\mathbb{Q}}\right) d\mathbb{Q} = KL(\mathbb{P}, \mathbb{Q})$$

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

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

• The joint convexity of KL, which we discuss below, is typically proved in a tedious way, relying on the rewriting of Exercise 1 and the joint convexity of
$$(x, y) \in [0, +\infty)^2 \longmapsto \left( \frac{x}{y} \ln \frac{x}{y} \right) y$$
We see it instead as a consequence of the data-processing inequality:

Corollary (joint convexity of KL):   For all probability distributions $\mathbb{P}_1, \mathbb{P}_2$ and $Q_1, Q_2$ over the same measurable space $(\Omega, \mathcal{F})$, and all $d \in (0, 1)$,

$$KL\left( (1-d)\, \mathbb{P}_1 + d\, \mathbb{P}_2, \ (1-d)\, Q_1 + d\, Q_2 \right) \leq (1-d)\, KL(\mathbb{P}_1, Q_1)$$
$$+ d\, KL(\mathbb{P}_2, Q_2)$$

Proof:   We augment $(\Omega, \mathcal{F})$ into $(\Omega \times \{1, 2\}, \mathcal{F}')$ where
$$\mathcal{F}' = \mathcal{F} \otimes \{ \emptyset, \{1\}, \{2\}, \{1, 2\} \}$$
We define the random pair $(X, J)$ by the projections
$$X : (\omega, j) \mapsto \omega \qquad \text{and} \qquad J : (\omega, j) \mapsto j$$

Let $\mathbb{P}$ be a probability measure on $(\Omega \times \{1, 2\}, \mathcal{F}')$ such that
$$\begin{cases} J \sim 1 + Ber(d) \\ X \mid J = j \sim \mathbb{P}_j \end{cases} \qquad \begin{array}{l}(\text{and a similar definition for } Q \\ \text{based on } Q_1, Q_2) \end{array}$$

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

Now, $\quad\quad\quad \mathbb{P}^X = (1-d)\,\mathbb{P}_1 + d\,\mathbb{P}_2$

$\quad\quad\quad\quad\quad Q^X = (1-d)\,Q_1 + d\,Q_2$

and (as we prove below) $\quad KL(\mathbb{P},Q) = (1-d)\,KL(\mathbb{P}_1, Q_1) + d\,KL(\mathbb{P}_2, Q_2)$

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

$\mathbb{P}_1 \ll Q_1 \quad$ and $\quad \mathbb{P}_2 \ll Q_2, \quad\quad$ so that $\quad\quad \mathbb{P} \ll Q \quad$ with

$$\frac{d\mathbb{P}}{dQ}(\omega, j) = \mathbb{1}_{\{j=1\}}\,\frac{d\mathbb{P}_1}{dQ_1}(\omega) + \mathbb{1}_{\{j=2\}}\,\frac{d\mathbb{P}_2}{dQ_2}(\omega)$$

This entails that (by Tonelli's theorem)

$$KL(\mathbb{P}, Q) = \int_{\Omega \times \{1,2\}} \left( \frac{d\mathbb{P}}{dQ}(\omega, j)\, \ln\frac{d\mathbb{P}}{dQ}(\omega, j) \right)\, dQ(\omega, j)$$

we integrate first over the $j$

$$= (1-d)\int_\Omega \left( \frac{d\mathbb{P}}{dQ}(\omega, 1)\, \ln\frac{d\mathbb{P}}{dQ}(\omega, 1) \right)\, dQ(\omega, 1)$$

$$+ \, d\int_\Omega \left( \frac{d\mathbb{P}}{dQ}(\omega, 2)\, \ln\frac{d\mathbb{P}}{dQ}(\omega, 2) \right)\, dQ(\omega, 2)$$

$$= \underbrace{(1-d)\int_\Omega \left( \frac{d\mathbb{P}_1}{dQ_1}(\omega)\, \ln\frac{d\mathbb{P}_1}{dQ_1}(\omega) \right)\, dQ_1(\omega)}_{= \, KL(\mathbb{P}_1, Q_1)} + \, d\,KL(\mathbb{P}_2, Q_2)$$

<u>KL for product measures.</u>    ($\leftrightarrow$ The independent case)

<u>Proposition:</u> Let $(\Omega, \mathcal{F})$ and $(\Omega', \mathcal{F}')$ be two measurable spaces,

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

$\mathbb{P}', Q'$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad$ over $(\Omega', \mathcal{F}')$

and denote by $\mathbb{P} \otimes \mathbb{P}'$ and $Q \otimes Q'$ the product distributions over

$(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$. Then:

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

**Proof:** It suffices to consider the case $P \ll P'$ and $Q \ll Q'$. Then $P \otimes P' \ll Q \otimes Q'$ with

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

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

Therefore, by Fubini-Tonelli (4. $x \mapsto x \ln x$ is bounded from below)

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

$$= \int_{\Omega} \left( \underbrace{\int_{\Omega} \left( \frac{dP}{dQ} \ln \frac{dP}{dQ} \right) dQ}_{= KL(P, Q)} \right) \underbrace{\frac{dP'}{dQ'} dQ'}_{dP'} + \underbrace{\text{similar term with } \ln \frac{dP'}{dQ'}}_{= KL(P', Q')}$$

$$= KL(P, Q)$$

**Consequence** (Garivier, Ménard, Stoltz, 2016):

_Data - processing inequality with expectations of random variables_

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

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

$$KL\left( Ber(E_P[X]), Ber(E_Q[X]) \right) \leq KL(P, Q)$$

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

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

$$KL\left( (P \otimes m)^{\mathbb{1}_E}, (Q \otimes m)^{\mathbb{1}_E} \right) \quad \leqslant \quad KL(P \otimes m, Q \otimes m)$$

$$\underbrace{\qquad}_{Ber(P \otimes m(E))} \quad \underbrace{\qquad}_{Ber(Q \otimes m(E))} \quad = \quad KL(P, Q) + KL(\eta, m)$$

$$\uparrow_{\text{if product distributions}} \quad = \quad KL(P, Q)$$

Thus: $\quad KL\left( Ber(P \otimes m(E)), Ber(Q \otimes m(E)) \right) \quad \leqslant \quad KL(P, Q)$

The proof is concluded by picking $E \in \mathcal{F} \otimes \mathcal{B}([q,1])$ such that

$$P \otimes m(E) = \mathbb{E}_P[x] \qquad \text{and} \qquad Q \otimes m(E) = \mathbb{E}_Q[x]$$

Namely, $\qquad E = \{ (\omega, x) \in \Omega \times [q,1] : \quad x \leqslant X(\omega) \}$

By Tonelli's theorem:

$$P \otimes m(E) = \int_\Omega \left( \int_{[q,1]} \mathbb{1}_{\{x \leqslant X(\omega)\}} \, dm(x) \right) dP(\omega)$$

$$= \int_\Omega X(\omega) \, dP(\omega) \quad = \quad \mathbb{E}_P[x]$$

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

<u>The chain rule</u> — 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.

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

A stochastic kernel $K$ is a mapping $(\Omega, \mathcal{F}) \to \mathcal{P}(\Omega', \mathcal{F}')$

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

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

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

$$\forall A \in \mathcal{F}, \qquad \forall B \in \mathcal{F}', \qquad K\mathbb{P}(A \times B) = \int_{\Omega} \mathbb{1}_A(\omega)\, K(\omega, B)\, d\mathbb{P}(\omega)$$

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

Theorem (chain rule for KL):

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

with $\quad$ (**) $g : (\omega, \omega') \mapsto \dfrac{dK(\omega, \cdot)}{dL(\omega, \cdot)}(\omega')$ being $\mathcal{F} \otimes \mathcal{F}'$-measurable,

Then

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

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

Question / Remark: $\qquad$ Not sure how needed my assumptions are!

$\qquad\qquad\qquad \hookrightarrow$ See Exercise 2 on the next page.

Proof: $\qquad \bullet \quad$ Since $g \ln g$ is lower bounded on $[0, +\infty)$, and in view of the measurability assumption (**), Tonelli's theorem wrt $LQ$ ensures that

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

is $\mathcal{F}$-measurable.

$\qquad\qquad \bullet \quad$ Given (*), we have $K\mathbb{P} \ll LQ$ if and only if $\mathbb{P} \ll Q$; we thus assume with no loss of generality that $K\mathbb{P} \ll LQ$ and $\mathbb{P} \ll Q$ (otherwise, both sides of the putative equality equal $+\infty$).

$\qquad\qquad \bullet \quad$ We write $\dfrac{d\mathbb{P}}{dQ} = f$, we then have

$$\frac{d\,K\mathbb{P}}{d\,L\mathbb{Q}}(\omega,\omega') = f(\omega)\,g(\omega,\omega')$$

as can be seen by going back to the definition of $L\mathbb{Q}$

And
$$KL(\,K\mathbb{P},\,L\mathbb{Q}) = \int_{\Omega\times\Omega'}\Big(f(\omega)\,g(\omega,\omega')\,\ln\big(f(\omega)g(\omega,\omega')\big)\Big)\,\underbrace{d\,L\mathbb{Q}(\omega,\omega')}_{L(\omega,d\omega')\,d\mathbb{Q}(\omega)}$$

by Tonelli's theorem

$$= \int_\Omega f(\omega)\,\ln f(\omega)\,\underbrace{\Big(\int_{\Omega'}g(\omega,\omega')\,L(\omega,d\omega')\Big)}_{\;=\,K(\omega,\Omega')\,=\,1}\,d\mathbb{Q}(\omega)$$

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

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

as announced.

## Exercise 2.

Try to weaken the assumptions of the chain-rule theorem.

→ Ideally, show that (*) and (**) can be relaxed into simply assuming (***) $\omega \mapsto KL(\,K(\omega,\cdot),\,L(\omega,\cdot))$ is $\mathcal{F}$-measurable

→ At least, I feel that one could prove

$$K\mathbb{P} \ll L\mathbb{Q} \quad\Longleftrightarrow\quad \begin{cases} \mathbb{P} \ll \mathbb{Q} \\[4pt] K(\omega,\cdot) \ll L(\omega,\cdot) \quad\text{for } \mathbb{P}\text{-almost all } \omega \end{cases}$$

In which case, we could assume with no loss of generality that (*) holds and the bi-measurability assumption (**) would be the only assumption.

Note: The only non-immediate implication is $K\mathbb{P} \ll L\mathbb{Q} \Rightarrow K(\omega,\cdot) \ll L(\omega,\cdot)$ for $\mathbb{P}$-almost all $\omega$.

<u>Lower bounds on the regret for stochastic bandits.</u>

Here is first a <u>summary</u> of the setting and context of stochastic bandits:

- K arms each indexed by $a = 1, 2, \ldots K$

- With each arm is associated a probability distribution $\nu_a \in \mathcal{D}$

- $\mathcal{D}$ 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 $\nu = (\nu_a)_{a \in \{1, \ldots K\}}$

- <u>Important quantities and notation:</u>

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

  $\mu^* = \max\limits_{a = 1 \ldots K} \mu_a$ is the largest expectation within $\nu$

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

  Arm $a$ is suboptimal if $\Delta_a > 0$

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

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

  1. The decision-maker picks $I_t \in \{1, \ldots K\}$ 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.

- Aim / regret: maximize $E\left[\sum\limits_{t=1}^{T} y_t\right]$

  which is equivalent to minimizing (controlling from above)

  $$R_T = T\mu^* - E\left[\sum\limits_{t=1}^{T} y_t\right]$$

- Rewriting by tower rule:

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

  where $N_a(T) = \sum\limits_{t=1}^{T} \mathbb{1}_{\{I_t = a\}}$ 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

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

$\underbrace{\qquad\qquad\qquad\qquad}$ history for the first $t$ rounds

$\underbrace{\qquad\qquad}$ arm picked at round $t+1$

- Strategies that are consistent wrt a model $\mathfrak{D}$:

If for all bandit problems $\vec{\nu} \in \mathfrak{D}^K$,

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

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

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

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

$$\forall \vec{\nu} \in \mathfrak{D}^K, \qquad \forall a \text{ s.t. } \Delta_a > 0,$$

$$\limsup_{T \to +\infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \leq C_a(\vec{\nu})$$

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

- Optimal (in some sense) such constant: $C_a(\vec{\nu}) = \dfrac{1}{K_{\inf}(\nu_a, \mu^*, \mathfrak{D})} = \dfrac{1}{K_{\inf}(\nu_a, \mu^*)}$

where $K_{\inf}(\nu_a, \mu^*, \mathfrak{D}) = K_{\inf}(\nu_a, \mu^*) = \inf\left\{ KL(\nu_a, \nu_a') ; \begin{array}{l} \nu_a' \in \mathfrak{D} \\ \mathbb{E}(\nu_a') > \mu^* \end{array} \right\}$

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

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

$\hookrightarrow$ The upper bound will come later.

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

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

For all strategies $\Psi$ consistent wrt $\mathfrak{D}$ (possibly randomized),

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

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

$$\liminf \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq 1 / K_{\inf}(\nu_a, \mu^*, \mathfrak{D})$$

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

For all (possibly randomized) strategies $\Psi$ consistent wrt $\mathcal{D}$,

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

$$\liminf_{T \to +\infty} \frac{\overline{R}_T}{\ln T} \geq \sum_{a : \Delta_a > 0} \frac{\Delta_a}{K_{\inf}(\nu_a, \mu^*, \mathcal{D})}.$$

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

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

Lemma (Fundamental inequality for stochastic bandits):

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

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

$$\sum_{a=1}^{K} E_{\vec{\nu}}[N_a(T)] \, KL(\nu_a, \nu_a') = KL(\mathbb{P}_{\vec{\nu}}^{H_T}, \mathbb{P}_{\vec{\nu}'}^{H_T})$$

$$\geq KL\left(Ber(E_{\vec{\nu}}[Z]), Ber(E_{\vec{\nu}'}[Z])\right)$$

Exercise 3:        Prove the theorem based on this lemma.

To that end:        fix $\mathcal{D}, \Psi, \vec{\nu} \in \mathcal{D}^K$, and $a$ s.t. $\Delta_a > 0$

Then build $\vec{\nu}' \in \mathcal{D}^K$ as $\begin{cases} \nu_j' = \nu_j & \forall j \neq a \\ \nu_a' \in \mathcal{D} & \text{with } E(\nu_a') > \mu^* \end{cases}$

$\vec{\nu}'$ is called an alternative problem (as opposed to $\vec{\nu}$: original bandit problem).

Note that $\begin{cases} a \text{ is suboptimal for } \vec{\nu} \\ a \text{ is the only optimal arm for } \vec{\nu}' \end{cases}$

and apply the lemma (fundamental inequality)        with $Z = N_a(T)/T$.