# Lecture #5: some properties of the KL

Last lecture, we proposed algorithms with (pseudo)regrets bounded as

$$R_T \leqslant c \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k} \qquad \text{(instance dependent regret)}$$

Is it possible to do better?

This lecture focuses on lower bounding the achievable regret by any algorithm

For that we consider a model where the rewards distributions belong to some **known** distribution set $\mathcal{D}$.

ie $\quad \forall k \in [K], \quad \nu_k \in \mathcal{D}$

  unknown $\qquad$ known

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

$$R_T \text{ is at best of order } \left( \sum_{k, \Delta_k > 0} \frac{\Delta_k}{K_{\inf}(\nu_k, \mu^*, \mathcal{D})} \right) \ln T$$

where $\qquad K_{\inf}(\nu_k, \mu^*, \mathcal{D}) = \inf \left\{ KL(\nu_k, \nu') \ \middle|\ \begin{array}{l} \nu' \in \mathcal{D} \\ \mathbb{E}[\nu'] > \mu^* \end{array} \right\}$

  Kullback-Leibler divergence

We will only prove the lower bound part

- **Case 1:** $\mathcal{D} = \left\{ \mathcal{N}(\mu, \sigma^2) \mid \mu \in \mathbb{R} \right\}$

then

$$K_{inf}\left(\nu_a, \mu^*, \mathcal{D}\right) = \frac{\Delta_a^2}{2\sigma^2}$$

Best possible regret of order $\quad 2\sigma^2 \sum_{a, \Delta_a > 0} \frac{\ln T}{\Delta_a}$

UCB has regret $\leq 32\sigma^2 \sum_{a, \Delta_a > 0} \frac{\ln T}{\Delta_a}$

$\hookrightarrow$ optimal up to constant factor

can be made optimal with finer version

- **Case 2:** $\mathcal{D} = \left\{ Ber(p) \mid p \in [0,1] \right\}$

then

$$K_{inf}\left(\nu_a, \mu^*, \mathcal{D}\right) = \mu_a \ln \frac{\mu_a}{\mu^*} + (1-\mu_a) \ln \frac{1-\mu_a}{1-\mu^*}$$

**But** before proving the lower bound, I guess that some reminder of basic and non-basic results about KL divergences would be needed!

# Definition

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

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

$\mathbb{Q}(A) = 0 \Rightarrow \mathbb{P}(A) = 0$

# Basic Facts

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

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

indeed, $\Psi$ is strictly convex. Jensen's inequality indicates that

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

equality if and only if $\frac{d\mathbb{P}}{d\mathbb{Q}}$ is $\mathbb{Q}$-almost surely constant, ie $\mathbb{P} = \mathbb{Q}$.

## A useful rewriting:

Assume $\mathbb{P} \ll \mathbb{Q}$ and let $\nu$ be any probability measure over $(\Omega, F)$ with $\mathbb{P} \ll \nu$, $\mathbb{Q} \ll \nu$. Denote $f = \frac{d\mathbb{P}}{d\nu}$, $g = \frac{d\mathbb{Q}}{d\nu}$,

then $KL(\mathbb{P}, \mathbb{Q}) = \int_{\Omega} \ln\left(\frac{f}{g}\right) f \, d\nu$.

see proof in exercise session 3.

useful when $\mathbb{P}$ and $\mathbb{Q}$ both admit densities over a classical reference
measure (eg Lebesgue).

## **Lemma** (data processing inequality)

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

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

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

Then $$KL(\mathbb{P}^X, \mathbb{Q}^X) \leqslant KL(\mathbb{P}, \mathbb{Q})$$

**Proof:** we can assume $P \ll Q$, since otherwise $KL(P,Q) = +\infty$ and it holds.

We show that we then have $P^X \ll Q^X$, with $\dfrac{dP^X}{dQ^X} = \underbrace{\mathbb{E}_Q\left[\dfrac{dP}{dQ}\Big| X = \cdot\right]}_{= \gamma}$

$\text{ie } \gamma(X) = \mathbb{E}_Q\left(\dfrac{dP}{dQ}\Big| X\right)$

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

$$P^X(B) = P(X \in B) = \int_\Omega \mathbb{1}_B(X) \frac{dP}{dQ} dQ \overset{\text{tower rule}}{=} \int_\Omega \mathbb{1}_B(X) \mathbb{E}_Q\left[\frac{dP}{dQ}\Big| X\right] dQ$$

$$= \int_\Omega \mathbb{1}_B(X) \gamma(X) dQ = \int_{\Omega'} \mathbb{1}_B \gamma \, dQ^X$$

by def of $Q^X$

therefore,

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

$$= \int_\Omega \left(\mathbb{E}_Q\left[\frac{dP}{dQ}\Big| X\right] \ln\left(\mathbb{E}_Q\left[\frac{dP}{dQ}\Big| X\right]\right)\right) dQ$$

$\psi$ is convex, conditional Jensen inequality

$$\leq \int_\Omega \mathbb{E}_Q\left[\frac{dP}{dQ} \ln \frac{dP}{dQ}\Big| X\right] dQ$$

tower rule

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

□

## References

- The proof above is due to Ali and Silvey ('66), 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, given below, is typically proved in a tedious way, relying on the joint convexity of $(x,y) \in \mathbb{R}_+^2 \longmapsto \left( x \ln \frac{x}{y} \right) y$. We may 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 $(\Omega, F)$ and all $\lambda \in (0,1)$:

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

**Proof:** We augment $(\Omega, F)$ into $(\Omega', F')$ where

$$\Omega' = \Omega \times \{1,2\}$$

$$F' = F \otimes \left\{ \varnothing, \{1\}, \{2\}, \{1,2\} \right\}$$

we define the random pair $(X, J)$ by the projections $X : \begin{matrix} \Omega \times \{1, 2\} \to \Omega \\ (\omega, j) \longmapsto \omega \end{matrix}$

and $J : \begin{matrix} \Omega \times \{1, 2\} \to \{1, 2\} \\ (\omega, j) \longmapsto j \end{matrix}$

Let $\mathbb{P}$ be a probability measure on $(\Omega', \mathcal{F}')$ such that:

$$\begin{cases} J \sim 1 + Ber(\lambda) \\ X | J = j \sim \mathbb{P}_j \end{cases}$$

(and a similar def for $\mathbb{Q}$ with $\mathbb{Q}_1, \mathbb{Q}_2$)

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

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

$\quad \mathbb{Q}^X = (1-\lambda) \mathbb{Q}_1 + \lambda \mathbb{Q}_2$

and <span style="color:blue">as we prove below</span> $\quad KL(\mathbb{P}, \mathbb{Q}) = (1-\lambda) KL(\mathbb{P}_1, \mathbb{Q}_1) + \lambda KL(\mathbb{P}_2, \mathbb{Q}_2) \quad$ so that the

result follows from the data processing inequality.

Indeed, we may assume with no loss of generality for $\lambda \in (0,1)$ that $\mathbb{P}_1 \ll \mathbb{Q}_1$, $\mathbb{P}_2 \ll \mathbb{Q}_2$, so that $\mathbb{P} \ll \mathbb{Q}$ with

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

This entails that:

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

$$= \int_{\Omega'} \left( \frac{d\mathbb{P}_1}{dQ_1}(\omega) \ln \frac{d\mathbb{P}_1}{dQ_1}(\omega) \right) \mathbb{1}_{\Omega \times \{1\}}(\omega, j) \, dQ(\omega, j)$$

$$+ \int_{\Omega'} \left( \frac{d\mathbb{P}_2}{dQ_2}(\omega) \ln \frac{d\mathbb{P}_2}{dQ_2}(\omega) \right) \mathbb{1}_{\Omega \times \{2\}}(\omega, j) \, dQ(\omega, j)$$

$$= \int_{\Omega} \left( \frac{d\mathbb{P}_1}{dQ_1}(\omega) \ln \frac{d\mathbb{P}_1}{dQ_1}(\omega) \right)(1-\lambda) \, dQ_1(\omega) + \cdots$$

$$= (1-\lambda) \, KL(\mathbb{P}_1, Q_1) + \lambda \, KL(\mathbb{P}_2, Q_2) \qquad \square.$$

**Proposition** (KL for product measures, independent case)

Let $(\Omega, F)$ and $(\Omega', F')$ be two measurable spaces

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

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

and denote by $\mathbb{P} \otimes \mathbb{P}'$ and $Q \otimes Q'$ the product distributions over $(\Omega \times \Omega', F \otimes F')$. Then

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

**Proof** we have $\mathbb{P} \otimes \mathbb{P}' \ll \mathbb{Q} \otimes \mathbb{Q}' \iff (\mathbb{P} \ll \mathbb{Q}$ and $\mathbb{P}' \ll \mathbb{Q}')$, so we can assume that all $\ll$ statements hold. Then

$$\frac{d(\mathbb{P} \otimes \mathbb{P}')}{d(\mathbb{Q} \otimes \mathbb{Q}')} = \frac{d\mathbb{P}}{d\mathbb{Q}} \quad \frac{d\mathbb{P}'}{d\mathbb{Q}'}$$

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

Therefore by Tonelli

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

if $\int\int f, g > 0$, then

$\int (f+g) \, d\mu = \int\int d\mu + \int g \, d\mu$

$$= \int_{\Omega'} \left( \int_{\Omega} \left( \frac{d\mathbb{P}}{d\mathbb{Q}} \ln \frac{d\mathbb{P}}{d\mathbb{Q}} \right) d\mathbb{Q} \frac{d\mathbb{P}'}{d\mathbb{Q}'} d\mathbb{Q}' \right. \qquad + \text{similar term with } \ln \frac{d\mathbb{P}'}{d\mathbb{Q}'}$$

$f = \frac{d\mathbb{P}'}{d\mathbb{Q}'}\left( \frac{d\mathbb{P}}{d\mathbb{Q}} \ln \frac{d\mathbb{P}}{d\mathbb{Q}} + \frac{1}{e} \right)$

$g = \frac{d\mathbb{P}}{d\mathbb{Q}}\left( \frac{d\mathbb{P}'}{d\mathbb{Q}'} \ln \frac{d\mathbb{P}'}{d\mathbb{Q}'} + \frac{1}{e} \right)$ here

$\underbrace{\phantom{KL(\mathbb{P},\mathbb{Q})}}_{KL(\mathbb{P},\mathbb{Q})}$  $d\mathbb{P}'$

$\underbrace{\phantom{KL(\mathbb{P},\mathbb{Q})}}_{KL(\mathbb{P},\mathbb{Q})}$  here we apply Tonelli again   $\underbrace{\phantom{KL(\mathbb{P}',\mathbb{Q}')}}_{KL(\mathbb{P}',\mathbb{Q}')}$

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

Data-processing inequality with expectations of random variables.

Let $X : (\Omega, F) \to ([0,1], B([0,1]))$ be any $[0,1]$-valued random variable

Then, denoting by $\mathbb{E}_{\mathbb{P}}[X]$ and $\mathbb{E}_{\mathbb{Q}}[X]$ the respective expectations of $X$ under $\mathbb{P}$ and $\mathbb{Q}$, we have:

$$\mathbb{E}_p[X] \ln \frac{\mathbb{E}_p[X]}{\mathbb{E}_Q[X]} + (1 - \mathbb{E}_p[X]) \ln \frac{1 \cdot \mathbb{E}_p[X]}{1 - \mathbb{E}_Q[X]} = KL\left(\text{Ber}(\mathbb{E}_p[X]), \text{Ber}(\mathbb{E}_Q[X])\right) \leqslant KL(P, Q)$$

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

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

$$KL\left((P \otimes \mu)^{\mathbb{1}_E}, (Q \otimes \mu)^{\mathbb{1}_E}\right) \leqslant KL(P \otimes \mu, Q \otimes \mu) = KL(P, Q) + KL(\mu, \mu)$$
$$= KL(P, Q).$$

$\underbrace{\qquad}_{\text{Ber}((P \otimes \mu)(E))} \quad \underbrace{\qquad}_{\text{Ber}((Q \otimes \mu)(E))}$

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

<span style="color:teal">Is it possible?</span>

Yes, taking $E = \left\{ (\omega, x) \in \Omega \times [0,1] : x \leqslant X(\omega) \right\} \in \mathcal{F} \otimes \mathcal{B}([0,1])$ as $X$ is measurable.

By Tonelli's theorem:

$$P \otimes \mu(E) = \int_\Omega \left( \int_{[0,1]} \mathbb{1}_{\{x \leqslant X(\omega)\}} \, d\mu(x) \right) dP(\omega) = \int_\Omega X(\omega) \, dP(\omega) \quad \text{and same for } Q.$$

$\square$

**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 $(\Omega, F)$ and $(\Omega', F')$ be two measurable spaces; we denote by $P(\Omega', F')$ the set of probability measures over $(\Omega', F')$.

A (regular) stochastic kernel $K$ is a mapping $(\Omega, F) \longrightarrow P(\Omega', F')$

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

such that $\forall B \in F', \omega \longmapsto K(\omega, B)$ is $F$-measurable

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

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

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

An extension of
Fubini (Tonelli) theorem

# Lemma

Let $\varphi : \Omega \times \Omega' \longrightarrow R$ be either $F \otimes F'$ measurable and $\geq 0$

or $KP$-integrable

Then $\quad w \longmapsto \int_{\Omega'} \varphi(w,w') K(w,dw')$ is $\mathcal{F}$-measurable and $\int_{\Omega\times\Omega'} \varphi \, dK\mathbb{P} = \int_\Omega \left( \int_{\Omega'} \varphi(w,w') K(w,dw') \right) d\mathbb{P}(w)$

**Proof:** (sketch) The result is true $\overset{\text{including measurability of } w \longmapsto \int \varphi(w,\cdot) K(w,d\circ) \text{ by regularity of } K}{\downarrow}$ for $\varphi = \mathbb{1}_{A\times B}$ by definition of $K\mathbb{P}$.

Extension to $\mathbb{1}_E$ for any $E \in \mathcal{F} \otimes \mathcal{F}'$ by an argument of $\sigma$-algebra contained /monotone class theorem, using monotone convergence (including the $w \longmapsto \int_{\Omega'} \cdots$ measurability)

Extension to $\begin{cases} \varphi \geq 0 & \text{by monotone convergence} \\ \varphi \in \mathbb{L}^1 \end{cases}$

$\overset{\text{actually with no loss of generality.}}{\downarrow}$

# Theorem (chain rule for KL): Assume $\mathbb{P} \ll \mathbb{Q}$

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

with $(**)$ the existence of a function $g: (w,w') \longmapsto \dfrac{dK(w,\cdot)}{dL(w,\cdot)}(w')$ being $\mathcal{F}\otimes\mathcal{F}'$-measurable, ↳ up to a $LQ$-null set

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

where $\quad w \longmapsto KL\big(K(w,\cdot), L(w,\cdot)\big)$ is indeed $\mathcal{F}$-measurable and $\geq 0$ so that the integral in the right-hand side is well defined.

# Remark:

1) the assumptions $(*)$ and $(**)$ will be satisfied for the

applications we have in mind.

2) They can be relaxed: - it suffices to assume that $\Omega'$ is a topological space with a countable base and $\mathcal{F}'$ is the Borel $\sigma$-algebra.

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

$$V = \bigcup_{i: O_i \subseteq V} O_i \quad , \text{ that is, as a countable union of elements of}$$

$(O_m)_{m \geq 1}$.

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

**Proof** ⊛ by bi-measurability of $g \ln g$, and since $g \ln g$ is lower bounded, 

<span style="color:blue">an immediate extension of</span> the previous lemma can be applied to get
$$w \longmapsto \int_{\Omega'} g(w,\cdot) \ln(g(w,\cdot)) \, L(w, d\cdot)$$
$$= KL(K(w,\cdot)) \, L(w,\cdot))$$

is $\mathcal{F}$-measurable and $\geq 0$

⊛ We assume $\mathbb{P} \ll \mathbb{Q}$, let $f = \frac{d\mathbb{P}}{d\mathbb{Q}}$. What can we say about $(w,w') \longmapsto f(w) g(w,w')$ ?

$$\int \mathbb{1}_{A \times B}(w,w') f(w) g(w,w') \, dL\mathbb{Q}(w,w') = \int_{\Omega} \left( \int_{\Omega'} \mathbb{1}_B(w') g(w,w') L(w, dw') \right) \mathbb{1}_A(w) f(w) \, d\mathbb{Q}(w)$$

<span style="color:blue">extension of Tonelli</span>

$$= \int_{\Omega'} \overbrace{1_B(\omega')\ K(\omega, d\omega')} = K(\omega, B)$$

$$= \int_{\Omega} \underbrace{1_A(\omega)\ K(\omega, B)}_{\text{F-measurable}}\ \underbrace{f(\omega)\ dQ(\omega)}_{dP(\omega)} = KP(A \times B) \qquad \text{by def } f\ KP$$

By Radon-Nikodym's Theorem: $\qquad \dfrac{dKP}{dLQ} = fg \qquad\qquad LQ\text{-as}$

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

  indeed $\quad LQ(A \times \Omega') = Q(A)$

  $\qquad\quad KP(A \times \Omega') = P(A)$.

- Therefore under (⊛), (⊛⊛), we have $\qquad KP \ll LQ \Longleftrightarrow P \ll Q$

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

$\Psi = fg \ln(fg)$ is lower bounded. The lemma (extension of Fubini-Tonelli extends to it):

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

<span style="color:blue">(again we can use the translation by $+\frac{1}{e}$ to justify this equality</span>

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

$$= \int_{\Omega} \left( KL\left( K(w,\cdot), L(w,\cdot) \right) + \ln(f(w)) \right) f(w) \, dQ(w)$$

$$= \int_{\Omega} KL\left( K(w,\cdot), L(w,\cdot) \right) \underbrace{f(w) \, dQ(w)}_{dP(w)} + \underbrace{\int f(w) \ln f(w) \, dQ(w)}_{KL(P,Q)}$$

□