Motivation: Histograms/Density → Comparing them, optimizing over them
→ ≠ Vectors, ≠ Functions ... WEAK OBJECTS

Measures: • On a set $X$ : $A \subset X \rightsquigarrow \mu(A) \in \mathbb{R}$ (vol, size, Mass) ①

$\left(\begin{array}{l} \text{Positive}: \mu(A) \geq 0 \qquad \text{Proba}: \mu(X) = 1 \\ \mu(A \cup B) = \mu(A) + \mu(B) \quad \text{if} \quad A \cap B = \emptyset \\ \qquad \hookrightarrow \text{should extend to contable union} \end{array}\right]$ 
→ Lebesgue "area"
→ Dirac
→ Density $\sum$ w$_i$ $\delta$ Diracs!

Radon measure: to speak about convergence, one needs a distance (of pair of measures)

① Needs all balls, $\mu(B) < +\infty$ (small enough) ② $\mu(A) = \sup \{ \mu(K) : \text{compact } K \subset A \}$
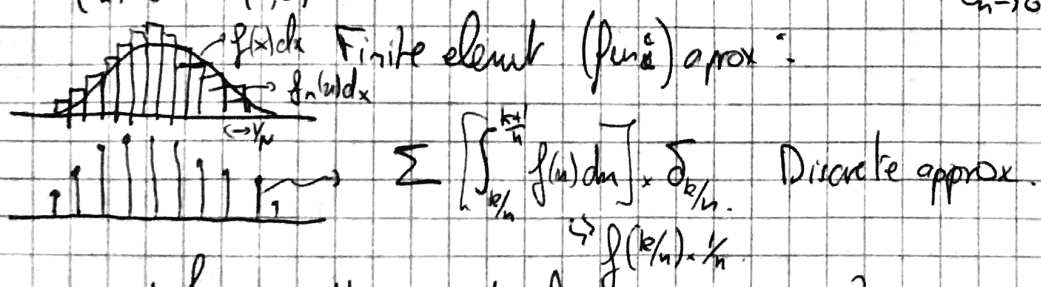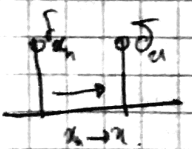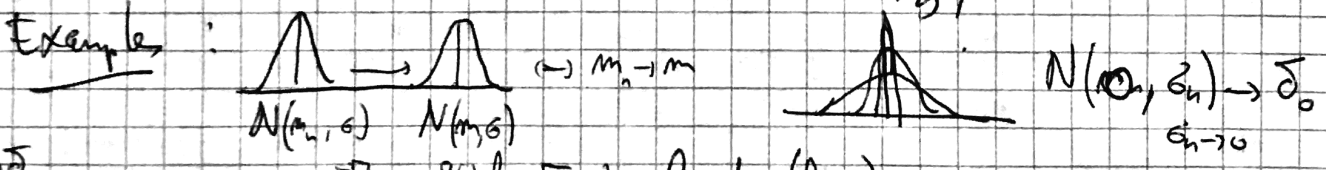
⇒ Allows to define $\int_A f \, d\mu$ for $f \in \mathcal{C}(X)$ continuous.

Radon Proba measures & Random variables: $Z: (\Omega, \mathbb{P}) \to X$ random var
associated measure $\mu(A) = \mathbb{P}_Z(A) \triangleq \mathbb{P}(Z \in A)$ proba
$Z$ is "push-forwarding" $\mathbb{P}$ to $\mathbb{P}_X$.

Convergence of random variables / Measures : convergence in law ⟷ weak* convergence

$\mu_n \xrightarrow{*} \mu \quad (\Longleftarrow) \quad \forall A, \mu_n(A) \to \mu(A)$
$\mathbb{P}_{Z_n} \xrightarrow{*} \mathbb{P}_Z \qquad (\Leftarrow) \quad \forall f \in \mathcal{C}(X), \int f \, d\mu_n \to \int f \, d\mu$

⚠ Weaker than the [STRONG] convergence of the density $\mu_n = g_n dx \quad g_n \xrightarrow{L^1} g \Rightarrow \mu_n \to \mu$
($\mu_n \xrightarrow{} \mu$ might not be have densities!) $\mu = g \, dx$ ⇏

Example:  $(\Leftarrow) m_n \to m$ 
$N(m_n, \sigma) \quad N(m, \sigma)$
$N(0, \sigma_n) \to \delta_0 \quad \sigma_n \to 0$

$\delta_{x_n} \quad \delta_x$ 
$x_n \to x$.

$f(x)dx$ Finite element (piecewise) approx
$f_n(x)dx$

$\sum \left[\int_{k/n}^{\frac{k+1}{n}} f(u) du\right] \cdot \delta_{k/n}$ Discrete approx.
$\hookrightarrow f(k/n) \cdot \frac{1}{n}$

Key question: Quantifying the speed of convergence? $D(\mu, \nu)$ "distance" like
First requirement: $\mu_n \to \mu \iff D(\mu_n, \mu) \xrightarrow{n} 0$
$(\Leftarrow) D$ defines the weak topology
WARNING: in $\infty$ dimension (and in particular if $D$ is not a norm)
$(D_1, D_2)$ same topology ⇎ equivalent $(\Leftarrow) \exists c, \frac{D_2}{c} \leq D_1 \leq c D_2$
ie convergence rates depend on the distance!

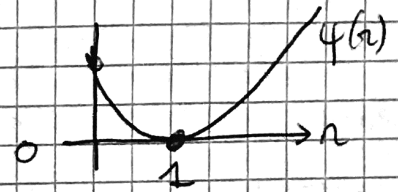# Statistical divergence

**① Comparing density:** if $\begin{cases} \mu = f\,dx \\ \nu = g\,dx \end{cases}$ then compute $D(\mu,\nu) = \left(\int |f(x)-g(x)|^p\,dx\right)^{1/p} = \|f-g\|_{L^p}$

Pbm: a strong assumption
• Not continuous with weak topology

$D(\mu + \varepsilon \delta, \nu)$ not defined

**② Comparing relative density:** ~~...~~ $\psi$-divergence / Csiszár divergence

Comparing $\dfrac{d\mu}{d\nu} = f$ (ie $d\mu(x) = f(x)\,d\nu$) with $1$

$$D(\mu|\nu) \triangleq \int_X \psi\left(\frac{d\mu}{d\nu}(u)\right) d\nu(u)$$
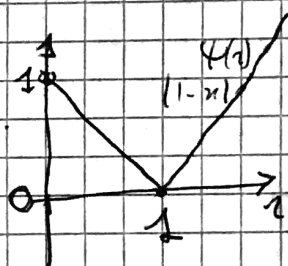


$\begin{cases} \psi \text{ convex} \\ \psi(1) = 0 \end{cases}$     $D(\mu|\nu) = 0 \Leftrightarrow \mu = \nu$
$D$ convex of $(\mu,\nu)$ !

examples: $\psi(x) = |x-1|$

$$D(\mu|\nu) = \int \left|\frac{d\mu}{d\nu} - 1\right| d\nu = \int \left|\frac{d\mu}{dx} - \frac{d\nu}{dx}\right| dx = \left\|\frac{d\mu}{dx} - \frac{d\nu}{dx}\right\|_{L^1}$$
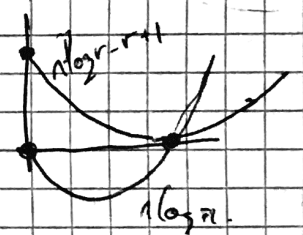
$TV(\mu-\nu) \quad \to$ it is a norm !!



example: $\psi(x) = x\log x$
  KL divergence

$$D(\mu|\nu) = \int \log\left(\frac{d\mu}{d\nu}\right) d\mu$$

$\left(\text{"generalized" KL}\right) \begin{cases} \psi(x) = x\log x - x + 1 \end{cases}$  $D(\mu|\nu) = \int \log\left(\frac{d\mu}{d\nu}\right) d\mu + \int(d\nu - d\mu)$

↳ useful to make it also a Bregman divergence
↳ same if $\int d\mu = \int d\nu$.



$x\log x - x + 1$

$x\log x$

**③ Hilbertian norm using lifting:** Fa simplicity, $X = \mathbb{R}^d$.

Convolution / Kernel density estimator: $\mu * h$ as a density   $\int h(x-y)\,d\mu(y)$
$\left(h(x) \text{ smooth } \int h = 1\right)$

ex  $\mu = \sum a_i \delta_{x_i} \longrightarrow \sum a_i h(x - x_i)$

$$D(\mu,\nu)^2 = \|h*\mu - h*\nu\|_{L^2}^2 = \int_{\mathbb{R}^d}\left[\int_{\mathbb{R}^d} h(x-y)\,d\xi(y)\right]^2 dx \triangleq \|\mu-\nu\|_k^2$$

$\to$ seems intractable even fa discrete $\mu$ and $\nu$

Max Mean (MMD) Discrepancy

$$= \int_{\mathbb{R}^d} \iint h(x-y)\,h(x-y')\,d\xi(y)\,d\xi(y')\,dx \quad \boxed{FUBINI}$$

$$= \iint K(y,y')\,d\xi(y)\,d\xi(y') \qquad K(y,y') \triangleq \int_{\mathbb{R}^d} h(x-y)h(x-y')\,dx$$
↳ Kernel

$\boxed{ex}$ $h(x) = e^{-\|x\|^2/2\sigma^2}$
$\rightsquigarrow K(y,y') = \exp(-\|y-y'\|^2/4\sigma^2)$

$\boxed{Ex}$ $K(y,y') = -\|y-y'\|$  also correspond to $h(x) = \frac{1}{\|x\|\cdots}$

So $D(\mu,\nu)^2 = \iint K(y,y')(d\mu(y)-d\nu(y))(d\mu(y')-d\nu(y'))$

$\boxed{Ok}$ $\mu = \sum a_k \delta_{x_k}$
$\nu = \sum b_k \delta_{y_k}$ $\longrightarrow$ $D(\mu,\nu)^2 = -2\sum_{i,j} k(x_i, y_j) a_i b_j$

$+ \sum_{i,i'} k(x_i, x_{i'}) a_i a_{i'}$

$+ \sum_{j,j'} k(y_j, y_{j'}) b_j b_{j'}$

$= E(k(X,X'))$
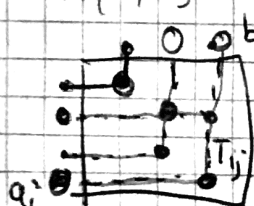$+ E(k(Y,Y'))$
$- 2 E(k(X,Y))$

$\boxed{Thm : \|\cdot\|_k \text{ metrizes weak convergence}}$

---

## Optimal Transport

Special discrete case : $\mu = \sum a_i \delta_{x_i}$  $\nu = \sum b_j \delta_{y_j}$  "Grain of masses"

$a_i, b_j \geq 0, \sum a_i = \sum b_j \ (=1)$.

Coupling $\Pi(a,b) = \{ P \in \mathbb{R}_+^{n \times m} : \forall_i \sum_j P_{ij} = a_i, \forall_j \sum_i P_{ij} = b_j \}$ cvx polytope

$P\mathbb{1} = a$  $P^T \mathbb{1} = b$

$\underline{OT :}$  $\mathcal{E}_c(\mu,\nu) = \min_{P \in \Pi(a,b)} \sum c_{ij} P_{ij}$ → Lin Prog / Combinatorial optim / Sinkhorn approx

$\underline{Wasserstein \ dist :}$ if $c_{ij} = d(x_i, y_j)^p$, $p \geq 1$, $W_p(\mu,\nu) \triangleq \mathcal{E}_c(\mu,\nu)^{1/p}$

$\boxed{Thm : W_p \text{ is a distance \& it metrizes weak convergence}}$

MMD : ⊕ Simple ⊕ Good Sample cxity ⊖ Reflects less the distance $\|\mu - \mu(\circ \tau)\| = O(\sqrt{\delta})$

OT ⊖ Cpx ⊖ Bad Sample cxity ⊕ More geometrical

$|D(\mu,\nu) - D(\hat\mu_n, \hat\nu_n)| \sim \frac{1}{n^q}$  MMD: $q = 1/2$  OT: $q = 1/d$