

Lecture 6:

Divergences between Probability measures

I Motivating problem: density fitting

• Fundamental problem: compare $\nu \in \mathcal{P}(\mathbb{R}^d)$ arising from measurements to a model which is a parameterized family of distributions $\{\mu_\theta; \theta \in \Theta\}$ where typically $\Theta \in \mathbb{R}^p$.

• A suitable parameter can be obtained by minimizing:

$$\min_{\theta \in \Theta} F(\theta) \text{ where } F(\theta) = D(\mu_\theta, \nu) \quad (*)$$

where $D: \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow [0, +\infty]$ is a divergence -

• In this lecture, by divergence we mean $\begin{cases} D(\mu, \nu) \geq 0 \\ D(\mu, \mu) = 0 \end{cases}$.

Example 1: One can choose $D(\mu, \nu) = W_p^p(\mu, \nu)$ for some $p \geq 1$.

When ν is an empirical measure, with $p=2$, the solution to (*) is called the Minimum Kantorovich Estimator -

Example 2: Let $x_1, \dots, x_n \in \mathbb{R}^d$ are independent samples from ν . When μ_θ has a density p_θ w.r.t a reference measure σ , the maximum likelihood estimator (MLE) is

$$\min_{\theta \in \Theta} -\frac{1}{n} \sum_{i=1}^n \log(p_\theta(x_i)) \quad (1)$$

This corresponds to an empirical version of solving (*) with $D(\mu, \nu) = \text{KL}(\mu, \nu)$ since (1) converges to $-\int \log(p_\theta(x)) d\nu(x) = \text{KL}(\mu_\theta, \nu) - \int \log\left(\frac{d\nu}{d\sigma}\right) d\nu$ (provided all the terms are finite) -

• Note that the MLE fails:

- when there is no natural reference measure σ
- when p_θ is difficult to compute
- when the objective F is too complicated to minimize.

Generative models

Generative models are when the parametric measure μ_θ is given by

$$\mu_\theta = (h_\theta)_\# \xi \quad \text{where } h_\theta: \mathbb{Z} \rightarrow \mathbb{R}^d$$

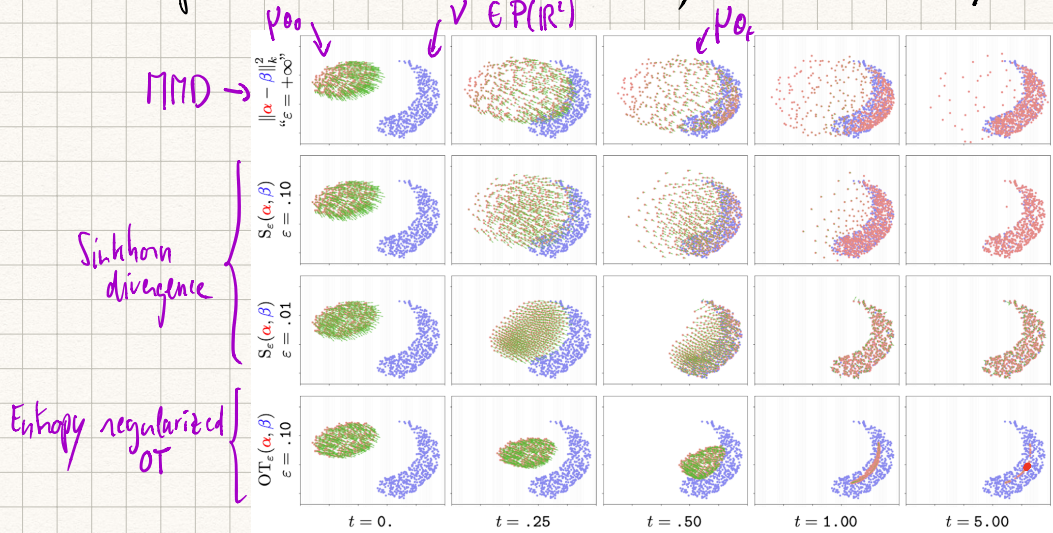
↑ latent space

and where $\xi \in \mathcal{P}(\mathbb{Z})$ is a reference measure. This leads to

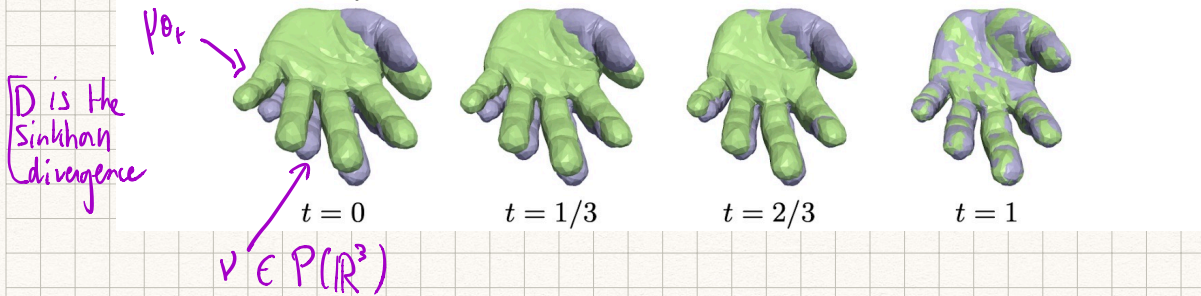
$$F(\theta) = D((h_\theta)_\# \xi, \nu).$$

The typical approach to "minimize" F is the gradient descent algorithm:

- initialize $\theta_0 \in \Theta$
 - for $t=1, 2, \dots$ let $\theta_{t+1} = \theta_t - \gamma \nabla F(\theta_t)$ where $\gamma > 0$ is a step-size
- formula?



Application to shape registration:



$(h_\theta)_\theta$ is a parameterized set of diffeomorphisms.

Let us give a formula for $\nabla F(\theta)$ under strong regularity assumptions.

Let us denote $E : \mu \mapsto D(\mu, \nu)$

Proposition. Assume that $E : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ is such that $\forall \mu \in \mathcal{P}_2(\mathbb{R}^d)$, there exists a function $E'(\mu) \in \mathcal{C}'(\mathbb{R}^d)$ with $\nabla E'(\mu)$ Lipschitz, and such that $\forall \nu \in \mathcal{P}(\mathbb{R}^d)$,

$$E(\nu) - E(\mu) = \int_{\mathbb{R}^d} E'(\mu) d(\mu - \nu) + o(W_2(\mu, \nu)).$$

Assume moreover that $h : \mathbb{R}^p \rightarrow L^2(\mathcal{S}; \mathbb{R}^d)$ is (Fréchet) differentiable, with partial derivatives at θ denoted $\partial_i h_\theta \in L^2(\mathcal{S}; \mathbb{R}^d)$. Then $F : \theta \mapsto E((h_\theta)_\# \mathcal{S})$ is differentiable with gradient, for $i=1, \dots, p$,

$$[\nabla F(\theta)]_i = \int_{\mathcal{S}} \nabla E'((h_\theta)_\# \mathcal{S})(h_\theta(z))^T \partial_i h_\theta(z) d\mathcal{S}(z).$$

← for digest in the practical session

Proof: First we study $G : f \mapsto E(f_\# \mathcal{S})$ and show that G is (Fréchet) differentiable with differential: $DG(f)(\delta f) = \int \nabla E'(f_\# \mathcal{S})(f(z))^T \delta f(z) d\mathcal{S}(z)$. Then the conclusion follows by the usual chain rule for Fréchet differentials.

For $f, \delta f \in L^2(\mathcal{S}, \mathbb{R}^d)$, we have that $W_2^2(f_\# \mathcal{S}, (f + \delta f)_\# \mathcal{S}) \ll \| \delta f \|_{L^2(\mathcal{S})}^2$ by taking $(f, f + \delta f)_\# \mathcal{S}$ as an admissible transport plan. Thus,

$$\begin{aligned} E((f + \delta f)_\# \mathcal{S}) - E(f_\# \mathcal{S}) &= \int_{\mathcal{S}} [E'(f_\# \mathcal{S})(f(z) + \delta f(z)) - E'(f_\# \mathcal{S})(f(z))] d\mathcal{S}(z) + o(\|\delta f\|) \\ &= \int_{\mathcal{S}} \nabla E'(f_\# \mathcal{S})(f(z))^T \delta f(z) d\mathcal{S}(z) + \underbrace{O(\text{Lip}(\nabla E'(f_\# \mathcal{S}) \|\delta f\|)}_{o(\|\delta f\|)} + o(\|\delta f\|) \end{aligned}$$

This shows $G(f + \delta f) - G(f) = DG(f)(\delta f) + o(\|\delta f\|)$. Hence the conclusion \square

Example: Show if $W : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is symmetric and differentiable with a Lipschitz gradient, then $E(\mu) := \int W(x, y) d\mu(x) d\mu(y)$ satisfies the assumptions above with $E'(\mu) : x \mapsto \int W(x, y) d\mu(y)$.

Now we will introduce various divergences and study: (i) the "divergence property"
 From now, X is a compact metric space. (ii) their weak continuity.

II Csizár divergences (a.k.a. f -divergences)

Definition. Let $f: \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. For $\mu, \nu \in \mathcal{P}(X)$, let $\mu = \left(\frac{d\mu}{d\nu}\right)\nu + \mu^\perp$ be the Lebesgue decomposition. We define

$$D_f(\mu, \nu) = \int f\left(\frac{d\mu}{d\nu}\right) d\nu + f'_\infty(1) \cdot \mu^\perp(X)$$

where $f'_\infty(x) := \lim_{t \rightarrow +\infty} f(tx)/t \in \mathbb{R} \cup \{+\infty\}$.

(\hookrightarrow recession/horizon of f)

Proposition: Let f be convex and such that $\min f = 0$ and $\arg\min f = \{1\}$.

Then $D_f(\mu, \nu) \geq 0$ with equality if and only if $\mu = \nu$.

Proof: If $\mu = \nu$ then $\frac{d\mu}{d\nu} = 1 \in L^1(\nu)$ and $\mu^\perp = 0$ so $D_f(\mu, \nu) = \int f(1) d\nu = 0$.
Conversely if $D_f(\mu, \nu) = 0$ then $\mu^\perp = 0$ (because $f'_\infty(1) \geq f(z) - f(1) > 0$) and $\frac{d\mu}{d\nu} = 1 \in L^1(\nu)$ so $\mu = \nu$.

Example (Kullback-Leibler divergence). Take $f(s) = \begin{cases} s \log s - s + 1 & \text{if } s > 0 \\ 1 & \text{if } s = 0 \\ +\infty & \text{if } s < 0 \end{cases}$

If $\mu \ll \nu$ then

$$D_f(\mu, \nu) = \int_X \left(\frac{d\mu}{d\nu} \log \left(\frac{d\mu}{d\nu} \right) - \frac{d\mu}{d\nu} + 1 \right) d\nu = \int_X \log \left(\frac{d\mu}{d\nu} \right) d\mu = \text{KL}(\mu, \nu),$$

and $D_f(\mu, \nu) = +\infty$ otherwise since $f'_\infty(1) = +\infty$.

Example (Total variation). Take $f(s) = \begin{cases} |s-1| & \text{if } s \geq 0 \\ +\infty & \text{otherwise} \end{cases}$

We have $f'_\infty(1) = 1$ thus

$$D_f(\mu, \nu) = \int_X \left(\left| \frac{d\mu}{d\nu} - 1 \right| d\nu + d\mu^\perp \right) \stackrel{(*)}{=} \int d|\mu - \nu| = \|\mu - \nu\|_{TV} =: \|\mu - \nu\|_{TV}$$

= sup_{f \in \mathcal{B}(X)} \int f d(\mu - \nu); \|f\|_\infty \leq 1

where $(*)$ comes from the fact that $\begin{cases} (\mu - \nu)_+ = \max\{0, \frac{d\mu}{d\nu} - 1\} \nu + \mu^\perp \\ (\mu - \nu)_- = \max\{0, 1 - \frac{d\mu}{d\nu}\} \nu \end{cases}$.

In the context of generative models, a drawback is that D_f is not weakly continuous in general: for instance $D_f(S_x, S_y) = \begin{cases} 0 & \text{if } x = y \\ f'_\infty(1) & \text{otherwise} \end{cases}$ is discontinuous in general.

Proposition. If $f: \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex, l.s.c and not identically $+\infty$, then $D_f(\mu, \nu)$ is (jointly) convex, weakly l.s.c and one has

$$D_f(\mu, \nu) = \sup_{\Psi \in \mathcal{C}(X)} \int \Psi d\mu + \int \Psi d\nu \text{ s.t. } \Psi(x) + f^*(\Psi(x)) \leq 0, \forall x \in X$$

where $f^*: s \mapsto \sup_{u \in \mathbb{R}} u \cdot s - f(u)$ is the convex conjugate of f .

Proof: see the lecture notes.

III Integral Probability Metrics (dual norms)

Definition: For a symmetric set B of measurable functions $X \rightarrow \mathbb{R}$ and $\alpha \in \mathcal{M}(X)$ a signed finite measure, let

$$\| \alpha \|_B := \sup_{f \in B} \int_X f(x) d\alpha(x)$$

For $\mu, \nu \in \mathcal{P}(X)$, with $\alpha = \mu - \nu$, we define

$$D_B(\mu, \nu) := \| \mu - \nu \|_B = \sup_{f \in B} \int f(x) d(\mu(x) - \nu(x))$$

This is called an "integral probability metric".

Proposition: If B is symmetric, bounded in sup-norm and contains 0, then $\| \cdot \|_B$ is a semi-norm on $\mathcal{M}(X)$, i.e. it is non negative, positively 1-homogeneous and subadditive.

Proof: left as an exercise.

Example 1: Total variation is recovered with $B = \{ f \in \mathcal{C}(X); \| f \|_\infty \leq 1 \}$.

Example 2: Wasserstein-1 (W_1) is recovered with $B = \{ f \in \mathcal{C}(X); \text{Lip}(f) \leq 1 \}$.

Example 3: The "flat norm" corresponds to

$$B = \{ f \in \mathcal{C}(X); \text{Lip}(f) \leq 1 \text{ and } \| f \|_\infty \leq 1 \}$$

To "metrize" the weak convergence, B should not be too large nor too small.

Proposition 3.5.

- (i) If $C(X) \subset \overline{\text{span}(B)}^{\|\cdot\|_\infty}$, i.e. the span of B is dense in $(C(X), \|\cdot\|_\infty)$ then
- $$\left\{ \begin{array}{l} \|\alpha_k - \alpha\|_B \rightarrow 0 \\ (\alpha_k) \text{ bounded for } \|\cdot\|_{TV} \end{array} \right. \text{ implies } \alpha_k \rightarrow \alpha$$
- (ii) If $B \subset C(X)$ is compact then
- $$\left\{ \begin{array}{l} \alpha_k \rightarrow \alpha \\ (\alpha_k) \text{ bounded for } \|\cdot\|_{TV} \end{array} \right. \text{ implies } \|\alpha_k - \alpha\|_B \rightarrow 0$$

Proof:

(i) If $\|\alpha_k - \alpha\|_B \rightarrow 0$, then $\forall f \in B$, since $\langle f, \alpha_k - \alpha \rangle \leq \|\alpha_k - \alpha\|_B$ so $\langle f, \alpha_k \rangle \rightarrow \langle f, \alpha \rangle$. By linearity, this extends to $\text{span}(B)$ and then to $\overline{\text{span}(B)}^{\|\cdot\|_\infty}$ since $|\langle f, \alpha_k \rangle - \langle f', \alpha_k \rangle| \leq \|f - f'\|_\infty \cdot \sup_k \|\alpha_k\|_{TV}$.

(ii) We assume that $\alpha_k \rightarrow \alpha$, consider a subsequence $(\alpha_{n_k})_k$ such that

$$\|\alpha_{n_k} - \alpha\|_B \rightarrow \limsup \|\alpha_k - \alpha\|_B$$

Since B is compact, let $f_{n_k} \in B$ achieve the supremum defining $\|\alpha_{n_k} - \alpha\|_B$.

We again extract a subsequence $(f_{n_k}) \xrightarrow{\|\cdot\|_\infty} f \in C(X)$. One has:

$$\|\alpha_{n_k} - \alpha\|_B = \langle \alpha_{n_k} - \alpha, f \rangle + \langle \alpha_{n_k}, f_{n_k} - f \rangle - \langle \alpha, f_{n_k} - f \rangle \rightarrow 0 \blacksquare$$

↳ This is a direct generalization of our proof of weak continuity of W_1 (in lecture 4).

IV Sinkhorn divergence

IV.1 Entropy Regularized optimal transport

Def (lecture 3). With $c \in \mathcal{C}(X \times X)$, let $\lambda \geq 0$ be the regularization, and

$$T_{c,\lambda}(\mu, \nu) := \min_{\gamma \in \Pi(\mu, \nu)} \int_{X \times X} c(x,y) d\gamma(x,y) + \lambda \text{KL}(\gamma, \mu \otimes \nu)$$

↑ differ by 1 from the one in lecture 3

Reminders:

(duality) $T_{c,\lambda}(\mu, \nu) = \sup_{\Psi, \Phi \in \mathcal{C}(X)} \int \Psi d\mu + \int \Phi d\nu + \lambda \left(1 - \int \int e^{(\Psi(x) + \Phi(y) - c(x,y))/\lambda} d\mu(x) d\nu(y) \right)$

(optimality cdt^o) - There exists maximizers $(\Psi_\lambda, \Phi_\lambda)$ and a unique minimizer γ_λ ,

linked by: $d\gamma_\lambda(x,y) = e^{(\Phi_\lambda(x) + \Psi_\lambda(y) - c(x,y))/\lambda} d\mu(x) d\nu(y)$

In particular, we have: $T_{c,\lambda}(\mu, \nu) = \int \Psi_\lambda d\mu + \int \Phi_\lambda d\nu$.

IV.2 Is $T_{c,\lambda}$ a suitable divergence?

Proposition: For $\mu, \nu \in \mathcal{P}(X)$, $c \in \mathcal{C}(X \times X)$, it holds:

$$T_{c,\lambda}(\mu, \nu) \rightarrow \begin{cases} T_c(\mu, \nu) = T_{c,0}(\mu, \nu) & \text{as } \lambda \rightarrow 0 \\ \int c(x,y) d\mu(x) d\nu(y) & \text{as } \lambda \rightarrow +\infty \end{cases}$$

see lecture 3

lecture notes.

Moreover, $\gamma_\lambda \rightarrow \mu \otimes \nu$ as $\lambda \rightarrow +\infty$.

Proof: see lecture notes.

$\gamma_\lambda^* = \int y d\nu(y)$ if $c(x,y) = \|y - x\|^2$

Corollary: let $\nu \in \mathcal{P}(X)$ be such that $\text{argmin}_{y \in X} \int c(x,y) d\nu(y)$ is a singleton $\{x^*\}$, and let

$$\mu_\lambda \in \text{argmin}_\mu T_{c,\lambda}(\mu, \nu).$$

Then as $\lambda \rightarrow +\infty$, one has $\mu_\lambda \rightarrow \delta_{x^*}$.

(proof see lecture notes).

IV.3 Debiased quantity: the Sinkhorn divergence

Thinking of $-T_{c,\lambda}$ as an "inner product" suggests to define

$$\underline{S_{c,\lambda}(\mu, \nu)} := T_{c,\lambda}(\mu, \nu) - \frac{1}{2} T_{c,\lambda}(\mu, \mu) - \frac{1}{2} T_{c,\lambda}(\nu, \nu)$$

Sinkhorn
divergence

Proposition (Interpolation properties). It holds, if $c(x, y) = \text{dist}(x, y)^p$ for $p \geq 1$,

$$S_{c,\lambda}(\mu, \nu) \rightarrow \begin{cases} T_c(\mu, \nu) & \text{as } \lambda \rightarrow 0 \\ \frac{1}{2} \|\mu - \nu\|_c & \text{as } \lambda \rightarrow \infty \end{cases}$$

where $\|\cdot\|_c$ is the kernel norm associated to $-c$.

(Proof is immediate from the previous proposition).

Proposition. • If $k(x, y) = e^{-c(x,y)/\lambda}$ is a p.o.d. kernel, then
 $S_\lambda(\mu, \nu) \geq 0$ with equality if $\mu = \nu$.

• If $e^{-c/\lambda}$ is furthermore a universal kernel, then
 $S_\lambda(\mu_n, \mu) \rightarrow 0$ if and only if $\mu_n \rightarrow \mu$.