


# A Primer on Optimal Transport


Probability distrib:  $(X, d)$  metric space eg  $X \subset \mathbb{R}^d$ ,  $d(x, y) = \|x - y\|$   
 For simplicity,  $X$  compact (otherwise needs moment conditions)

Radon measures:  $\alpha \in \mathcal{M}(X)$ :  $A \subset X \mapsto \alpha(A) \in \mathbb{R}$  Borel measure with regularity condition  
 $\mathcal{M}(X) = \mathcal{C}(X)^*$  dual of cont. funct:  $\int_A d\alpha(x)$

ie  $f \in \mathcal{C}(X) \mapsto \int_X f(x) d\alpha(x)$  is linear continuous

ex:  $\alpha = \sum_{i=1}^n a_i \delta_{x_i}$  discrete measure  $\int f d\alpha = \sum f(x_i) a_i$  

2 "ways" to discretize  $\swarrow$  weight  $(a_i)$ : EMER  
 $\searrow$  points  $(x_i)$ : LAGRANGE

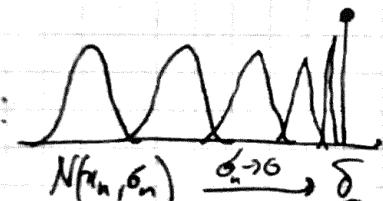
ex:  $d\alpha(x) = p(x) dx \Leftrightarrow \frac{d\alpha}{dx} = p$    $\int f d\alpha = \int f \cdot p dx$

Proba measures:  $[\alpha \geq 0 \text{ and } \int d\alpha = 1] \Leftrightarrow [\text{Law of random variable } X]$

not:  $\alpha \in \mathcal{M}_+^1(X)$   $\mathbb{E}_X(f(X)) = \int f(x) d\alpha(x)$

Weak\* topology:  $\alpha_n \xrightarrow{*} \alpha \Leftrightarrow \forall f, \int f d\alpha_n \rightarrow \int f d\alpha$

(w in law)  $X_n \rightrightarrows X \Leftrightarrow \forall f, \mathbb{E}(f(X_n)) \rightarrow \mathbb{E}(f(X))$   
 $\Leftrightarrow \forall A \subset X, \mathbb{P}(X_n \in A) \rightarrow \mathbb{P}(X \in A)$

ex:   $\delta_{x_n} \xrightarrow{*} \delta_x$

Comparing distrib:  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ , convex,  $\varphi(1) = 0$

$D_\varphi(\alpha | \beta) \triangleq \int_X \varphi\left(\frac{d\alpha}{d\beta}(x)\right) d\beta(x) + \alpha^\perp(X) \varphi_\infty$  Cizon  $\varphi$ -div

Prop  $D_\varphi(\cdot | \cdot) \geq 0$  is jointly cvx,  $D_\varphi(\alpha | \beta) = 0 \Leftrightarrow \alpha = \beta$

ex:  $\varphi(s) = \log(s)$   $D_\varphi = \text{KL} = \text{rel. entropy}$

$\varphi(s) = |1-s|$   $D_\varphi(\alpha, \beta) = \left\| \frac{d\alpha}{d\beta} - 1 \right\|_{L^1}$  TV

$\varphi(s) = \sqrt{1-s^2}$   $D_\varphi(\alpha, \beta) = \left\| \sqrt{\frac{d\alpha}{d\beta}} - \sqrt{\frac{d\beta}{d\alpha}} \right\|_{L^2}^2$  Hellinger<sup>2</sup>

$\varphi(s) = s \log(s) - (s-1) \log(s-1)$   $D_\varphi(\alpha, \beta) = \text{KL}(\alpha | \frac{\alpha+\beta}{2}) + \text{KL}(\beta | \frac{\alpha+\beta}{2})$  (Hellinger is a dist)  
 $(JS^{1/2} \text{ is a distance})$

②



Pbm: Dp doesn't metrize cv in low:  $\alpha + \alpha$ ,  $\| \alpha - \beta \|_{TV} = 2$   $\alpha = \delta_{x_1}$ ,  $\beta = \delta_{x_2}$   
 $\leadsto$  One wants  $D(\alpha|\beta)$  such that  $\alpha_n \rightarrow \alpha \iff D(\alpha_n|\alpha) \rightarrow 0$

Kernel norms: aka MMD (maximum Mean Discrepancy)

$k(x,y)$  positive kernels:  $\forall (x_i)_{i=1}^n, (k(x_i, x_j))_{i,j=1}^n > 0$

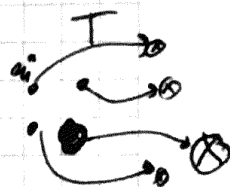
then:  $\| \xi \|_k^2 \triangleq \iint_{X^2} k(x,y) d\xi(x) d\xi(y)$  is a norm

ex:  $\xi = \sum_{i=1}^n a_i \delta_{x_i}$   $\| \xi \|_k^2 = \sum_{i,j} a_i a_j k(x_i, x_j) \leadsto n^2$  ops

Prop: if  $k$  universal (ie  $\text{span } k(x, \cdot)$  dense in  $\mathcal{C}(X)$ ) then  $\| \alpha - \beta \|_k^2$  metrizes weak convergence.

Ex: Gaussian, ED

Monge OT: Push-Forward:  $T: X \rightarrow X$   $T_\# : \delta_x \rightarrow \delta_{T(x)}$



By linearity:  $T_\# : \sum a_i \delta_{x_i} \rightarrow \sum a_i \delta_{T(x_i)}$

For measure  $\beta = T_\# \alpha$  defined by  $\forall f \in \mathcal{C}(X), \int f(y) d\beta(y) = \int f(T(x)) d\alpha(x)$

For random variables  $Y \sim \beta = T_\# \alpha$  is defined as  $Y = T(X)$ .

For density  $\leadsto$  painful.

Monge Pbm:  $\min_{T: X \rightarrow X} \left\{ \int_X d(x, T(x))^p d\alpha(x) : T_\# \alpha = \beta \right\}$

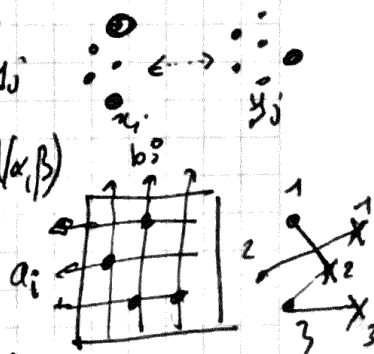
Pbm: In general, constraints are empty:  $\text{---} \rightarrow \text{---}$

Non convex intractable pbm  $\rightarrow$  convex relaxation.

Discrete Kantorovich pbm:  $\alpha = \sum_i a_i \delta_{x_i} \leftrightarrow \beta = \sum_j b_j \delta_{y_j}$

Coupling:  $\{ P \in \mathbb{R}_+^{m \times m} : \sum_j P_{ij} = a_i, \sum_i P_{ij} = b_j \} = \mathcal{U}(\alpha, \beta)$

$W_p(\alpha, \beta)^p \triangleq \min \{ \sum P_{ij} d(x_i, y_j)^p : P \in \mathcal{U}(\alpha, \beta) \}$



$\leadsto$  extends to arbitrary measures: replace  $P \leadsto \pi \in \Pi(X^2)$

$\mathcal{U}(\alpha, \beta) = \{ \pi \in \Pi(X^2) : \pi_1 = \alpha, \pi_2 = \beta \}$   $\Sigma \leadsto \int_X$

Examples: (1D)  $W_p(\sum_{i=1}^n \delta_{x_i}, \sum_{j=1}^m \delta_{y_j}) = \| \text{sort}(x) - \text{sort}(y) \|$ ,  $W_2(N(m_1, \sigma_1^2), N(m_2, \sigma_2^2)) = \|m_1 - m_2\| + \|\sigma_1^2 - \sigma_2^2\|$  (3)

Monge  $\leftrightarrow$  Kantorovich: Discrete:  $\sum_{i=1}^n \delta_{x_i} = \alpha, \sum_{j=1}^m \delta_{y_j} = \beta \Rightarrow \exists P \text{ admissible Permut. matrix (Birkhoff-Von Neumann)}$

Continuous (GRENIER):  $P$  optimal supported on  $x \mapsto T(x)$ ,  $T(x) = \nabla \psi(x)$ ,  $\psi$  convex  $\rightarrow$  MONGE AMPERE

Prop:  $W_p$  is a distance on  $\mathcal{M}_p^+(\mathbb{X})$  which metrizes weak<sup>\*</sup> conv

ex:  $W_p(\delta_x, \delta_y) = d(x, y) \rightarrow$  not possible with MMD

Entropic Regulariz<sup>o</sup>:  $\min_{P \in \mathcal{U}(a, b)} \underbrace{\sum_{i,j} P_{ij} d(x_i, y_j)^p}_{\substack{\text{arbitrary} \\ \text{measures}}} + \epsilon \underbrace{KL(P | \alpha \otimes \beta)}_{\sum_{i,j} P_{ij} \log(\frac{P_{ij}}{\alpha_i \beta_j})}$

SCHRODINGER PBM  $\rightarrow W_p^P(a, b) \triangleq \min_{\pi \in \Pi_+(X)} \int d\pi + \epsilon KL(\pi | \alpha \otimes \beta) : \pi_1 = \alpha, \pi_2 = \beta$

Prop:  $P$  optimal  $\Leftrightarrow \exists (u, v) \in \mathbb{R}^n \times \mathbb{R}^m$   $P_{ij} = K_{ij} u_i v_j, a, b; K = e^{-d^p/\epsilon}$   
unique and  $P1 = a, P^T 1 = b$

Cor: one needs to find  $(u, v)$  st.  $\begin{cases} P1 = (u \otimes a) \otimes (K \times (v \otimes b)) = a \\ P^T 1 = (v \otimes b) \otimes (K^T \times (u \otimes a)) = b \end{cases}$

Sinkhorn Algorithm:  $u \leftarrow \frac{1}{K(v \otimes b)} \leftrightarrow v \leftarrow \frac{1}{K(u \otimes a)}$  complexity:  $O(n^2)$  per iter<sup>o</sup>

THM/Sinkhorn (7): converges!

log-domain:  $(u, v) = (e^{f/\epsilon}, e^{g/\epsilon})$ ,  $(f, g)$  are "dual variables"

$$f_i \leftarrow -\epsilon \log \left[ \sum_j \exp \left( -\frac{d(x_i, y_j)^p - g_j}{\epsilon} \right) b_j \right] = \text{SoftMin}_\epsilon^b [C(x_i, \cdot) - g(\cdot)]$$

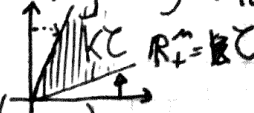
LSE  $\rightarrow$  stabilize

$$g_j \leftarrow \text{SoftMax}_\epsilon^a [C(\cdot, y_j) - f(\cdot)] \quad C = d^p$$

Approximation of the c-transforms:  $f^c \triangleq \min_{x \in \mathbb{X}} c(x, \cdot) - f(x)$

linear convergence rates: Variation norm:  $\|f\|_V = \inf_{c \in \mathbb{R}} \|f + c\|_\infty$   $\rightarrow$  defined up to constant

Hilbert metric:  $d_H(u, u') = \|\log(u) - \log(u')\|_\infty$

Contract<sup>o</sup>:  $\exists \eta \in ]0, 1[$ ,  $d_H(Ku, Ku') \leq \eta d_H(u, u')$  

$\rightarrow$  linear rate of Sinkh.  $d_H(1/u, 1/u') = d_H(u, u')$

④

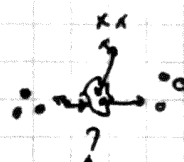
Complexity :  $O(n^2)$  per iteration,  $O(n)$  in some case (eg Gaussian convol)  
Matrix/vector mult<sup>n</sup>  $\xrightarrow{\text{parallelize}}$  Matrix/Matrix mult<sup>n</sup>

Bonus: Sinkhorn divergences:  $W_p^p(\alpha, \beta) - \frac{1}{2}W_p(\alpha, \alpha) - \frac{1}{2}W_p(\beta, \beta) \geq 0 \xrightarrow{\varepsilon \rightarrow \infty} \|\alpha - \beta\|_{-d}^2$

Bonus: Sample complexity of MMD vs Sinkhorn !!!

Barycenters :  $(\alpha, \beta) \rightsquigarrow W_p^p(\alpha, \beta)$  convex (look at the dual plan).

[Carlier/Aguade 2011]  $\min_{\alpha} \sum_k \lambda_k W_p^p(\alpha, \beta_k)$



If  $\beta_k = \delta_{x_k} \rightarrow$  usual def<sup>n</sup>  $\min_{\alpha} \sum_k \lambda_k \|\alpha - x_k\|^p$

If  $\beta_k$  discrete  $\Rightarrow$   $\alpha^*$  is discrete  $\beta_k$  Gaussian  $\Rightarrow \alpha^*$  Gaussian

Entropic regul<sup>n</sup> :  $\min_{\alpha} \sum_k \lambda_k W_p^p(\alpha, \beta_k) = \min_{\alpha, (\pi_k)_k} \sum_k \lambda_k KL(\pi_k | K_{\alpha, \beta_k})$   
 $\rightsquigarrow$  Sinkhorn-like iterations  $(\pi_k)_1 = \alpha$   $(\pi_k)_2 = \beta_k$  Rel. mass.

Gradient-flows :  $\min_{x \in \mathbb{R}^d} f(x) \rightsquigarrow \frac{dx_t}{dt} = - \nabla f(x_t)$   
with respect to  $\langle \cdot, \cdot \rangle_{\mathbb{R}^d}$

Explicit:  $x_{k+1} = x_k - \tau \nabla f(x_k)$  ( $t \approx \tau k$ )

Implicit:  $x_{k+1} + \tau \nabla f(x_{k+1}) = x_k \rightsquigarrow x_{k+1} = (Id + \tau \nabla f)^{-1}(x_k)$

"Proximal point algo" ie:  $x_{k+1} = \text{Argmin}_{x} \frac{1}{2} \|x - x_k\|^2 + \tau f(x)$   $\text{Prox}_{\tau f}(x_k)$

Advantages : no need for  $f$  to be smooth, no constraint on  $f$ ,

Disadvantage : need to solve a sub-problem often intractable

Fundamental remark : derivative free  $\rightarrow$  works over a metric-space !

$\alpha_{k+1} \triangleq \text{Argmin}_{\alpha} d(\alpha_k, \alpha)^p + \tau f(\alpha)$

Take  $d(x, y) = \|x - y\|$ ,  $p=2 \rightarrow$  Implicit Euler.

Take  $d(\alpha, \beta) = W_2(\alpha, \beta)$ ,  $p=2 \rightarrow$  Jordan-Kinderlehrer-Otto Wans. Flow

Continuous flow : let  $\tau \rightarrow 0$  then  $\alpha_k \xrightarrow{\tau k \rightarrow t} \alpha(t)$  for some class of function  $f$

"measure valued" derivative :  $f(\alpha + \epsilon \xi) = f(\alpha) + \epsilon \int f'(\alpha) d\xi + o(\epsilon)$   $f'(\alpha) \in \mathcal{C}(X)$

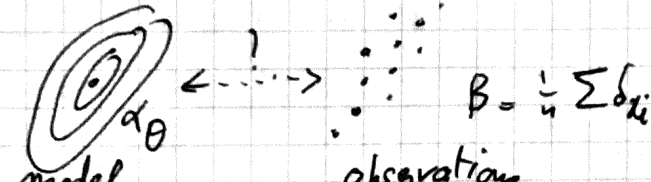
Formally,  $\alpha(t)$  solves :  $\frac{\partial \alpha}{\partial t} = -\text{div} [f'(\alpha) \cdot \nabla (f'(\alpha))]$

Example :  $f(\alpha) = \int u d\alpha \leadsto$  Advection  $\frac{\partial \alpha}{\partial t} = \text{div}(\alpha \nabla u)$

$f(\alpha) = \int \log(\frac{d\alpha}{dx}) d\alpha \leadsto$  Heat  $\Delta \alpha = \frac{\partial \alpha}{\partial t}$

$f(\alpha) = \int h(\frac{d\alpha}{dx}) d\alpha \leadsto$  non-linear PDE, eg Porous Medium  $\Delta \alpha$

$f(\alpha) = \iint k(x,y) d\alpha(x) d\alpha(y) \leadsto$  pairwise inter., Keller-Segel, etc

Density-Fitting :   $\alpha_\theta$  model  $\beta = \frac{1}{n} \sum \delta_{x_i}$  observations

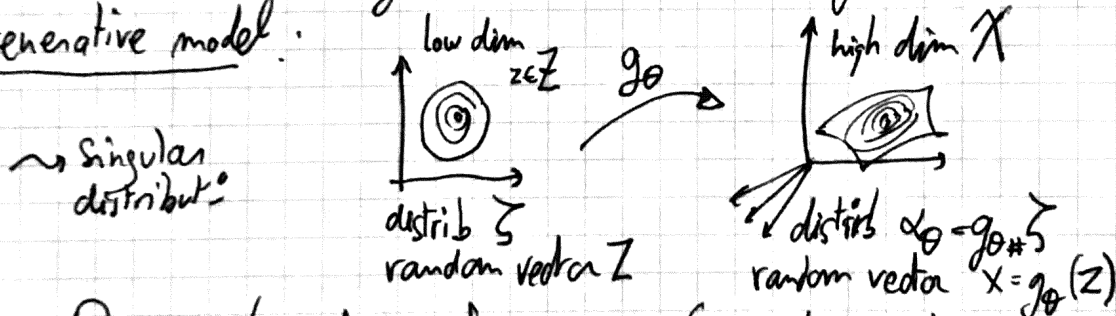
Maximum-likelihood :  $d\alpha_\theta = p_\theta dx \leadsto$  fixed for all  $\theta$  (only its support matters)

$\min_\theta \frac{1}{n} \sum \log(p_\theta(x_i)) \xrightarrow{n \rightarrow \infty} KL(\beta_n | \alpha)$  Ex: Gaussian estimat., Mixture (EM)

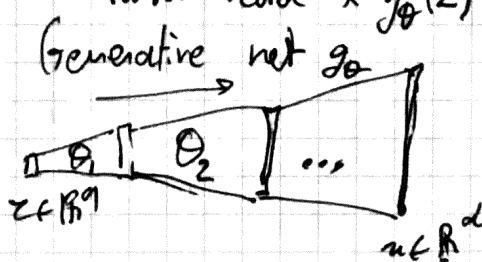
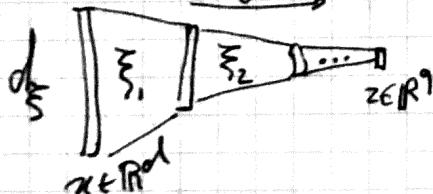
Prob: if support  $b_\theta$  "moves" +  $p_\theta$  might be intractable...

Minimum divergence estimator :  $\min_\theta D(\alpha_\theta | \beta)$

Generative model :



Discriminatory deep-net



$\varphi$ -divergence : not ok, not weak\* continuous }  $\rightarrow$  use duality!  
OT / MMD : not tractable, need to know a metric

Other approach :  $OT \rightarrow$  Sinkhorn } + metric learning (of Aude Generag)  
MMD

⑥

$$\alpha \rightarrow D(\alpha, \beta) \text{ convex, } D(\alpha, \beta) = \sup_{f: X \rightarrow \mathbb{R}} \int f(x) d\alpha(x) - D^*(f, \beta)$$

$\int f(x) d\alpha(x)$   
Legendre transform w.r.t  $\alpha$

Example:  $\varphi$ -div:  $D_\varphi^*(f, \beta) = \int \varphi^*(f(x)) d\beta(x)$

TV:  $\varphi(x) = |1-x| \rightarrow \varphi^*(t) = L_{[-\frac{1}{2}, \frac{1}{2}]}(t)$

KL:  $\varphi(x) = x \log x \rightarrow \varphi^*(t) = \exp(t-1)$

JS:  $\varphi(x) = x \log x - (x+1) \log(x+1) \rightarrow \varphi^*(t) = -\log(1-e^t) + L_{\mathbb{R}^+}(t)$

Example: MMD:  $D(\alpha, \beta) = \frac{1}{2} \|\alpha - \beta\|_K^2$ ,  $D^*(f, \beta) = \int f d\beta - \frac{1}{2} \|f\|_{RKHS}^2$

Example:  $W_1$ :  $D(\alpha, \beta) = W_1(\alpha, \beta)$ ,  $D^*(f, \beta) = \int f d\beta - L_{Lips}(f)$   
 $\hookrightarrow \|f\|_{Lips} \leq 1$

Min-Max fit:  $\min_{\theta} \max_f \int_X f d\alpha_\theta - D^*(f, \beta)$   
 $\int_Z f(g_\theta(z)) d\zeta(z)$

Generative Adversarial Networks: restrict  $f = d_\zeta$  to be discr. net

[Goodfellow et al.]

$$\min_{\theta} \max_{\zeta} \int_Z d_\zeta(g_\theta(z)) d\zeta(z) - D^*(d_\zeta, \beta)$$

$\xrightarrow{\text{for } \varphi\text{-div.}} = \frac{1}{n} \sum_i \varphi^*(d_\zeta(x_i))$

Concl:  $\rightarrow$  hard to solve but impressive result

$\rightarrow$  Restricting  $f = d_\zeta$ : not anymore a true divergence

$\rightarrow$  Unclear which  $D$  is best (is OT really useful?)