Generative Models and Optimal Transport

Marco Cuturi



Joint work / work in progress with G. Peyré, A. Genevay (ENS), F. Bach (INRIA), G. Montavon, K-R Müller (TU Berlin)

Statistics 0.1 : Density Fitting



Statistics 0.1 : Density Fitting



Density Fitting



Density Fitting



Density Fitting



ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By R. A. Fisher, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma



 $\max_{\boldsymbol{\theta}\in\Theta}\frac{\mathbf{1}}{N}\sum_{i=1}\log \boldsymbol{p}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})$

 $\nu_{\rm data}$

 $\nu_{\rm data}$

PA.

ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By R. A. Fisher, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma



 $\max_{\boldsymbol{\theta}\in\Theta}\frac{1}{N}\sum_{i}\log \boldsymbol{p}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})$ i=1

 $\log 0 = -\infty$ $p_{\theta}(x_i) \text{ must be } > 0$





In higher dimensional spaces...



In higher dimensional spaces...



In higher dimensional spaces...



















Push-forward: $\forall B \subset \Omega, \mathbf{f}_{\sharp}\boldsymbol{\mu}(B) := \boldsymbol{\mu}(\mathbf{f}^{-1}(B))$







Difference between fitting a push forward measure $f_{\theta \sharp} \mu vs.$ a density p_{θ} ?











• Formulation as adversarial problem [GPM...'14]

 $\min_{\boldsymbol{\theta} \in \Theta} \max_{\text{classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$

• Use a **richer metric** Δ for probability measures, able to handle measures with non-overlapping supports:

$$\min_{\boldsymbol{\theta}\in\Theta} \Delta(\boldsymbol{\nu}_{data}, \boldsymbol{p}_{\boldsymbol{\theta}}), \quad \min_{\boldsymbol{\theta}\in\Theta} \operatorname{KL}(\boldsymbol{\nu}_{data} \| \boldsymbol{p}_{\boldsymbol{\theta}})$$

Minimum Δ Estimation

The Annals of Statistics 1980, Vol. 8, No. 3, 457-487

MINIMU 1 CHI-SQUARE, NOT MAXIMUM LIKELIHOOD!

By JOSEPH BERKSON Mayo Clinic, Rochester, Minnesota



COMPUTATIONAL STATISTICS & DATA ANALYSIS

VIER Computational Statistics & Data Analysis 29 (1998) 81-103



Minimur Hellinger listance estimation for Poisson mixtures

Dimitris Karlis, Evdokia Xekalaki* Department of Statistics. Athens University of Economics and Business, 76 Patizsian Str., 104 34 Athens, Greece



Available online at www.sciencedirect.com

SCIENCE DIRECT.



Statistics & Probability Letters 76 (2006) 1298-1302

www.elsevier.com/locate/stapro

On minimum Kantorovich listance estimators

Federico Bassetti^a, Antonella Bodini^b, Eugenio Regazzini^{a,*}

Minimum Kantorovich Estimation

• Use optimal transport theory, namely *Wasserstein* distances to define discrepancy Δ .

$$\min_{\boldsymbol{\theta}\in\Theta} W(\boldsymbol{\nu}_{\text{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

• Optimal transport? fertile field in mathematics.



A geometric toolbox to compare **probability measures** supported on a metric space.



A geometric toolbox to compare probability measures supported on a <u>metric space</u>.



A powerful **geometric toolbox** to compare **probability measures**.





A powerful **geometric toolbox** to compare **probability measures**.



[SDPC..'15]

A powerful **geometric toolbox** to compare **probability measures**.

Wasserstein Barycenters [**Agueh'11**]

[SDPC..'15]
666. Mémoires de l'Académie Royale

MÉMOIRE

SUR LA

THÉORIE DES DÉBLAIS

ET DES REMBLAIS.

Par M. MONGE.

L'orsqu'on doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.





















 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T : \Omega \to \Omega$ $\inf_{T_{\sharp} \mu = \nu} \int_{\Omega} c(x, T(x)) \mu(dx)$



 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T : \Omega \to \Omega$ [Brenier'87] If $\Omega = \mathbb{R}^d, c = \| \cdot - \cdot \|^2$, μ, ν a.c., then $T = \nabla u, u$ convex.



Monge's Problem

 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T : \Omega \to \Omega$ $\inf_{T_{\sharp} \mu = \nu} \int_{\Omega} c(x, T(x)) \mu(dx)$



Monge's Problem

 Ω a probability space, $\boldsymbol{c}: \Omega \times \Omega \to \mathbb{R}$. $\boldsymbol{\mu}, \boldsymbol{\nu}$ two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $T: \Omega \to \Omega$ $\inf_{\mathbf{T}_{\sharp}\boldsymbol{\mu}=\boldsymbol{\nu}}\int_{\Omega}\boldsymbol{c}(x,\boldsymbol{T}(x))\boldsymbol{\mu}(dx)$ $oldsymbol{\delta_x}$ 26

[Kantorovich'42] Relaxation

• Instead of maps $T : \Omega \to \Omega$, consider probabilistic maps, i.e. couplings $P \in \mathcal{P}(\Omega \times \Omega)$:

$$\Pi(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega, \\ \boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \\ \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B}) \}$$

[Kantorovich'42] Relaxation

 $\Pi(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega,$ $\boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B})\}$



[Kantorovich'42] Relaxation

 $\Pi(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathcal{P}(\Omega \times \Omega) | \forall \boldsymbol{A}, \boldsymbol{B} \subset \Omega,$ $\boldsymbol{P}(\boldsymbol{A} \times \Omega) = \boldsymbol{\mu}(\boldsymbol{A}), \boldsymbol{P}(\Omega \times \boldsymbol{B}) = \boldsymbol{\nu}(\boldsymbol{B}) \}$



Wasserstein Distances

Def. For $p \ge 1$, the *p*-Wasserstein distance between μ, ν in $\mathcal{P}(\Omega)$, defined by a metric D on Ω ,

 $W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \iint \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})^p \boldsymbol{P}(d\boldsymbol{x}, d\boldsymbol{y}).$ PRIMAL

Wasserstein Distances

Def. For $p \ge 1$, the *p*-Wasserstein distance between μ, ν in $\mathcal{P}(\Omega)$, defined by a metric D on Ω ,

$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \iint \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})^p \boldsymbol{P}(d\boldsymbol{x}, d\boldsymbol{y}).$

THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES

BY FRANK L. HITCHCOCK

1. Statement of the problem. When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary from city to city.

МАТЕМАТИЧЕСНИЕ МЕТОДЫ

A.B. HANTOPOBHY

О РГАНИЗАЦИИ И ПЛАНИРОВАНИЯ ПРОИЗВОДСТВА

Wasserstein Distances

Def. For $p \ge 1$, the *p*-Wasserstein distance between μ, ν in $\mathcal{P}(\Omega)$, defined by a metric D on Ω ,

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \inf_{\boldsymbol{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \iint \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})^p \boldsymbol{P}(dx, dy).$$

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\substack{\boldsymbol{\varphi} \in L_1(\boldsymbol{\mu}), \boldsymbol{\psi} \in L_1(\boldsymbol{\nu})\\ \boldsymbol{\varphi}(x) + \boldsymbol{\psi}(y) \leq \boldsymbol{D}^p(x, y)}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

W is versatile



W is versatile



Minimum Kantorovich Estimators

 $\min_{\boldsymbol{\theta}\in\Theta} W(\boldsymbol{\nu}_{\mathrm{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$

- [Bassetti'06] 1st reference discussing this approach.
- [MMC'16] use regularization in a finite setting.
- [ACB'17] (WGAN) [BJGR'17] (Wasserstein ABC).
- **Hot topics**: <u>approximate</u> & <u>differentiate</u> W efficiently.
- Today: ideas from our recent preprint [GPC'17]







Consider
$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$.
 $M_{XY} \stackrel{\text{def}}{=} [D(x_i, y_j)^p]_{ij}$
 $U(a, b) \stackrel{\text{def}}{=} \{ P \in \mathbb{R}^{n \times m} | P \mathbf{1}_m = a, P^T \mathbf{1}_n = b \}$
 $\stackrel{\mathbf{y}_1 \qquad \cdots \qquad \mathbf{y}_m}{:} \stackrel{\mathbf{b}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{b}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_1 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf{i}_2 \qquad \cdots \qquad \mathbf{i}_m}{:} \stackrel{\mathbf$

Consider
$$\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\boldsymbol{\nu} = \sum_{j=1}^{m} b_j \delta_{y_j}$.
 $M_{\boldsymbol{X}\boldsymbol{Y}} \stackrel{\text{def}}{=} [D(\boldsymbol{x}_i, \boldsymbol{y}_j)^p]_{ij}$
 $U(\boldsymbol{a}, \boldsymbol{b}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathbb{R}^{n \times m}_+ | \boldsymbol{P} \boldsymbol{1}_m = \boldsymbol{a}, \boldsymbol{P}^T \boldsymbol{1}_n = \boldsymbol{b} \}$

Def. Optimal Transport Problem $W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$









Note: flow/PDE formulations [**Beckman'61**]/[**Benamou'98**] can be used for *p*=1/*p*=2 for a sparse-graph metric/Euclidean metric.








Discrete OT Problem



Discrete OT Problem



Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$

$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij}(\log P_{ij})$$

Note: Unique optimal solution because of strong concavity of Entropy

Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



Note: Unique optimal solution because of strong concavity of Entropy

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle P, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} \log P_{ij} + \alpha^T (P\mathbf{1} - \mathbf{a}) + \beta^T (P^T \mathbf{1} - \mathbf{b})$$
$$\partial L/\partial P_{ij} = M_{ij} + \gamma (\log P_{ij} + 1) + \alpha_i + \beta_j$$
$$(\partial L/\partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma} + \frac{1}{2}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma} + \frac{1}{2}} = \mathbf{u_i} K_{ij} \mathbf{v_j}$$

Fast & Scalable Algorithm

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

- [Sinkhorn'64] fixed-point iterations for $(\boldsymbol{u}, \boldsymbol{v})$ $\boldsymbol{u} \leftarrow \boldsymbol{a}/K\boldsymbol{v}, \quad \boldsymbol{v} \leftarrow \boldsymbol{b}/K^T\boldsymbol{u}$
- O(nm) complexity, GPGPU parallel [C'13].
- $O(n^{d+1})$ if $\Omega = \{1, \ldots, n\}^d$ and D^p separable. [S.C.'15]

Sinkhorn Divergence

Def. For
$$\gamma > 0$$
, let $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P_{\gamma}}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$

Prop.
$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\mu}) > 0$$

Def. Normalized Sinkhorn Divergence

$$\bar{W}_{\gamma}(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} W_{\gamma}(\boldsymbol{\mu},\boldsymbol{\nu}) - \frac{1}{2} \left(W_{\gamma}(\boldsymbol{\mu},\boldsymbol{\mu}) + W_{\gamma}(\boldsymbol{\nu},\boldsymbol{\nu}) \right)$$

Prop. If
$$p = 1$$
, $\overline{W}_{\gamma}(\mu, \nu) \xrightarrow[\gamma \to \infty]{} ED(\mu, \nu)$

Algorithmic Formulation

Def. For $L \geq 1$, define $W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P}_L, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle,$ where $P_L \stackrel{\text{def}}{=} \operatorname{diag}(\boldsymbol{u}_L) K \operatorname{diag}(\boldsymbol{v}_L)$, $\boldsymbol{v_0} = \boldsymbol{1}_m; l \ge 0, \boldsymbol{u_l} \stackrel{\text{def}}{=} \boldsymbol{a} / K \boldsymbol{v_l}, \boldsymbol{v_{l+1}} \stackrel{\text{def}}{=} \boldsymbol{b} / K^T \boldsymbol{u_l}.$ **Prop.** $\frac{\partial W_L}{\partial X}, \frac{\partial W_L}{\partial a}$ can be computed recursively, in O(L) kernel $K \times$ vector products.

Proposal: Autodiff OT using Sinkhorn

Approximate W loss by the transport cost \overline{W}_L after L Sinkhorn iterations.



[GPC'17]

Example: MNIST, Learning f_{θ}



Example: Generation of Images



MMD-GAN

 $\tau = 1000$

 $\tau = 10$

- CIFAR 10 images
- In these examples the cost function is also learned adversarially, as a NN mapping onto feature vectors.

Concluding Remarks

- *Regularized* OT is much faster than OT.
- *Regularized* OT can interpolate between W and the *MMD / Energy distance* metrics.
- The solution of *regularized OT* is *"auto-differentiable"*.
- Many open problems remain!

