Information Geometry of Entropy-Regularized Optimal Transport

Ryo Karakida (National Institute of AIST, Japan)

Shun-ichi Amari (RIKEN CBS, Japan),

Masafumi Oizumi (Univ. of Tokyo, Japan),

Marco Cuturi (Google Brain & ENSAE, France)

ICIAM2019 minisymposium "Distance Metrics and Mass Transfer Between High Dimensional Point Clouds" July 17th , 2019 "Information geometry of Wasserstein divergence" Ryo Karakida, Shun-ichi Amari Geometric Science of Information (GSI 2017)

"Information Geometry Connecting Wasserstein Distance and Kullback-Leibler Divergence via the Entropy-Relaxed Transportation Problem" Shun-ichi Amari, Ryo Karakida, Masafumi Oizumi Information Geometry (2018)

"Iformation Geometry for Regularized Optimal Transport and Barycenters of Patterns" Shun-ichi Amari, Ryo Karakida, Masafumi Oizumi, Marco Cuturi Neural Computation, (2019)







Outline

- Background
 - Entropy-regularized Optimal Transport (OT)
- Information geometry of entropy-regularized OT
 - Optimal transportation plan as exponential family
 - Dually flat structure
- Alternative divergence to entropy-regularized cost
 - An information geometric viewpoint
 - Barycenter of patterns

Optimal Transport Problem in S_{n-1}

(a.k.a. Hitchcock Problem)

Probability simplex: •

$$S_{n-1} = \left\{ p \in \mathbb{R}^n \mid \sum_i p_i = 1, \quad p_i \ge 0 \right\}$$

$$p_1 \qquad p_2 \qquad q_1$$

$$p_2 \qquad q_2$$

$$p_1 \qquad p_2 \qquad q_2$$

$$p_2 \qquad q_2$$

$$p_1 \qquad p_2 \qquad q_2$$

$$p_1 \qquad q_3$$

$$p_1 \qquad q_4$$

$$p_4$$

Optimal *P* is solvable by linear programming

 q_1

 $\mathbf{Q} = \mathbf{q}_i$

 M_{ij}

Optimal Transport Problem in S_{n-1}

Wasserstein distance:

$$W(\boldsymbol{p}, \boldsymbol{q}) = \min_{P \in \Pi(\boldsymbol{p}, \boldsymbol{q})} \sum_{i,j} M_{ij} P_{ij}$$
$$\boldsymbol{p}, \boldsymbol{q} \in S_{n-1}$$

Difficulties

- Computationally demanding $O(n^3 \ln n)$
- Solution is not necessarily unique
- Undifferentiable with ${m p}$ and ${m q}$

Entropy-regularized OT [M. Cuturi, NIPS 2013]:

$$W(\boldsymbol{p}, \boldsymbol{q}) = \min_{P \in \Pi(\boldsymbol{p}, \boldsymbol{q})} \sum_{i,j} M_{ij} P_{ij} + \lambda \sum_{ij} P_{ij} \ln P_{ij} \quad (\lambda > 0)$$

Advantages • Fast O (n^2) and GPU-friendly

- Unique solution
- Differentiable

Entropy-R Optimal Transport in S_{n-1}

• Barycenter of images q_i (i = 1, 2, ..., N) $\min_{p} \sum_{i=1}^{N} W(p, q_i)$



[Cuturi&Doucet, ICML 2014]

Spatially close pixels tend to take similar values

KL divergence



Wasserstein



Application



Imaging [Solomon+, SIGGRAPH 2015]

Color Grading [Bonneel+, SIGGRAPH 2016]



Dictionary learning [Schmitz+ 2018], generative models [Genevay+ 2017], ...

Kullback-Leibler (KL) divergence

Invariant under reversible transformations of random variables



Information Geometry [Amari 1985]: Fisher information metric (Riemannian metric), dual affine connections

Geometry of probabilistic distributions

Kullback-Leibler (KL) divergence

Invariant under reversible transformations of random variables

$$\int dp(x) \ln \frac{p(x)}{q(x)}$$

Information Geometry [Amari 1985]



Wasserstein distance

Reflects the ground metric supporting probability measures

$$\left(\inf_{P\in\Pi(\mu,\nu)}\int d(x,y)^p dP(x,y)\right)^{1/p}$$



Is **Information Geometry** not related?

Geometry of probabilistic distributions

Kullback-Leibler (KL) divergence

Invariant under reversible transformations of random variables

$$\sum_{i} p_i \ln \frac{p_i}{q_i}$$

Information Geometry [Amari 1985]

Wasserstein distance

Reflects the ground metric supporting probability measures





We can introduce Information Geometry of OT plans

Outline

- Background
 - Entropy-regularized Optimal Transport (OT)
- Information geometry of entropy-regularized OT
 - Optimal transportation plan as exponential family
 - Dually flat structure
- Alternative divergence to entropy-regularized cost
 - An information geometric viewpoint
 - Barycenter of patterns

Entropy-regularized OT in S_{n-1}

Probability simplex:
$$S_{n-1} = \left\{ \boldsymbol{p} \in \mathbb{R}^n \mid \sum_i p_i = 1, \quad p_i \ge 0 \right\}$$

Marginal distributions: $\boldsymbol{p}, \boldsymbol{q} \in S_{n-1}$

Metric between p_i and $q_j : M_{ij}$

Entropy-regularized OT [M. Cuturi, NIPS 2013]: $\varphi_{\lambda}(\boldsymbol{p}, \boldsymbol{q}) = \min_{\mathbf{P}} \frac{1}{1+\lambda} \langle \mathbf{M}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}) \quad (\lambda > 0)$

$$\langle \mathbf{M}, \mathbf{P} \rangle = \sum_{ij} M_{ij} P_{ij}$$
, $H(\mathbf{P}) = -\sum P_{ij} \log P_{ij}$

How to solve: Method of Lagrange multiplier

$$L_{\lambda}(\mathbf{P}) = \frac{1}{1+\lambda} \langle \mathbf{M}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}) - \sum_{i,j} (\alpha_i + \beta_j) P_{ij} \quad \text{Multipliers: } \alpha_i, \beta_j$$

Entropy-regularized OT: Unique Solution

$$L_{\lambda}(\mathbf{P}) = \frac{1}{1+\lambda} \langle \mathbf{M}, \mathbf{P} \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}) - \sum_{i,j} (\alpha_i + \beta_j) P_{ij} \qquad \text{Multipliers: } \alpha_i, \beta_j$$

[M. Cuturi, NIPS 2013] revealed

Entropy-regularized OT has a unique optimal solution P_{λ}^* $P_{\lambda ij}^* = ca_i b_j K_{ij}$ with $K_{ij} = \exp\left(-\frac{M_{ij}}{\lambda}\right)$ $a_i = \exp\left(\frac{1+\lambda}{\lambda}\alpha_i\right), \quad b_j = \exp\left(\frac{1+\lambda}{\lambda}\beta_j\right), \quad c = \frac{1}{\sum a_i b_j K_{ij}}$

• Multipliers (α, β) are determined by iterative computations

Sinkhorn algorithm: $oldsymbol{a} \leftarrow oldsymbol{p}./K^Toldsymbol{b}, \ oldsymbol{b} \leftarrow oldsymbol{q}./Koldsymbol{a}$

Entropy-regularized OT: Convexity

• Cost function φ_{λ} $\varphi_{\lambda}(\boldsymbol{p},\boldsymbol{q}) = \frac{1}{1+\lambda} \langle \mathbf{M}, \mathbf{P}_{\lambda}^* \rangle - \frac{\lambda}{1+\lambda} H(\mathbf{P}_{\lambda}^*)$

Lemma. Entropy-regularized cost $\varphi_{\lambda}(\boldsymbol{p}, \boldsymbol{q})$ convex over $(\boldsymbol{p}, \boldsymbol{q})$

Proof sketch : Wasserstein distance (linear in P) + Entropy (convex over P)

• Set of OT plans is an **Exponential Family**

$$P_{\lambda}^{*}(x;\boldsymbol{\theta}) = \exp\left\{\sum_{i,j} \theta^{ij} \delta_{ij}(x) - \psi_{\lambda}\right\} \qquad \frac{\theta^{ij} = \frac{1+\lambda}{\lambda}(\alpha_{i}+\beta_{j}) - \frac{M_{ij}}{\lambda}}{\delta_{ij}(x) = 1 \text{ when } x = (i,j), \\ 0 \text{ otherwise}}$$

- Normalization factor ψ_{λ} is convex over (α , β)
- 2(n-1)-dimensional manifold; $oldsymbol{lpha},oldsymbol{eta} \in \mathbb{R}^{n-1}$

 $\boldsymbol{\eta} = (\boldsymbol{p}, \boldsymbol{q})^T \qquad \boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta})^T$

Theorem. Normalization factor ψ_{λ} and cost φ_{λ} are both convex and connected by the Legendre Transformation;

$$\psi_{\lambda}(\boldsymbol{\theta}) + \varphi_{\lambda}(\boldsymbol{\eta}) = \boldsymbol{\theta} \cdot \boldsymbol{\eta} , \quad \boldsymbol{\theta} = \nabla_{\boldsymbol{\eta}} \varphi_{\lambda}(\boldsymbol{\eta}), \quad \boldsymbol{\eta} = \nabla_{\boldsymbol{\theta}} \psi_{\lambda}(\boldsymbol{\theta})$$

$$\psi_{\lambda}(\boldsymbol{\theta}) + \varphi_{\lambda}(\boldsymbol{\eta}) = \boldsymbol{\theta} \cdot \boldsymbol{\eta} \quad \boldsymbol{\theta} = \nabla_{\boldsymbol{\eta}} \varphi_{\lambda}(\boldsymbol{\eta}), \quad \boldsymbol{\eta} = \nabla_{\boldsymbol{\theta}} \psi_{\lambda}(\boldsymbol{\theta})$$

• The set of P_{λ}^* is a 2(n - 1)-dimensional **dually flat manifold**



Riemannian metric of dually flat manifold :

$$\mathbf{G}_{\lambda} = \nabla_{\boldsymbol{\eta}} \nabla_{\boldsymbol{\eta}} \varphi_{\lambda}(\boldsymbol{\eta}), \quad \mathbf{G}_{\lambda}^{-1} = \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \psi_{\lambda}(\boldsymbol{\theta})$$

Fisher information matrix (in the θ -coordinates) is explicitly given by $\mathbf{G}_{\lambda}^{-1} = \begin{bmatrix} \underline{p_i \delta_{ij} - p_i p_j | P_{ij} - p_i q_j} \\ P_{ij} - p_i q_j | q_i \delta_{ij} - q_i q_j \end{bmatrix}$

- Useful in estimation; $\min_{oldsymbol{q}} arphi_{\lambda}(oldsymbol{p},oldsymbol{q})$

Dual affine connections : By using a cubic tensor $T_{ijk} = \partial_i \partial_j \partial_k \psi_\lambda$, $\Gamma_{ijk} = [ij;k] - \frac{1}{2}T_{ijk}$, $\Gamma^*_{ijk} = [ij;k] + \frac{1}{2}T_{ijk}$

[*ij*; *k*]: Levi-Civita connection

Sinkhorn algorithm = iterations of e-projection

[Amari, Karakida & Oizumi, Information Geometry (2018)]

Sinkhorn algorithm:

• e-projection to $M(\boldsymbol{p}, \cdot)$

$$\begin{bmatrix} \boldsymbol{a} \leftarrow \boldsymbol{p}./K^T \boldsymbol{b} \\ \boldsymbol{b} \leftarrow \boldsymbol{q}./K \boldsymbol{a} \end{bmatrix} P^*_{\lambda i j} = c a_i b_j K_{i j}$$

 $M(m{p}, \cdot)~$: Subspace conditioned by a fixed $m{p}$ $M(\cdot, m{q})~$: Subspace conditioned by a fixed $m{q}$

$$T_{A}.\mathbf{P} = (a_i P_{ij})$$
 , $a_i = rac{p_i}{\sum_j P_{ij}}$

 $KL[T_A.\mathbf{P}:\mathbf{P}] + KL[\mathbf{P}^*:T_A.\mathbf{P}] = KL[\mathbf{P}^*:\mathbf{P}]$ $KL[T_A.\mathbf{P}:\mathbf{P}^*] \le KL[\mathbf{P}:\mathbf{P}^*]$

- e-projection to $M(\cdot, q)$

$$T_{\cdot B} \mathbf{P} = (b_j P_{ij})$$
 , $b_j = \frac{q_j}{\sum_i P_{ij}}$



Outline

- Background
 - Entropy-regularized Optimal Transport (OT)
- Information geometry of entropy-regularized OT
 - Optimal transportation plan as exponential family
 - Dually flat structure
- Alternative divergence to entropy-regularized OT[AKOC'19]
 - An information geometric viewpoint
 - Barycenter of patterns

Entropy regularized cost:

$$C_{\lambda}(\boldsymbol{p},\boldsymbol{q}) = \min_{P \in \Pi(\boldsymbol{p},\boldsymbol{q})} \langle \mathbf{M}, \mathbf{P} \rangle - \lambda H(\mathbf{P}) = (1+\lambda)\varphi_{\lambda}(\boldsymbol{p},\boldsymbol{q})$$

- It does **not** satisfy a criterion of divergence (distance); $\exists q \ C_{\lambda}(p,p) > C_{\lambda}(p,q)$
- Minimum is given by $q^* = ilde{K}_\lambda p$

Conditional metric

$$\tilde{\mathbf{K}}_{\lambda} = \left[\frac{K_{ji,\lambda}}{K_{j\cdot}}\right]_{ij} \quad K_{j\cdot} = \exp\left(-\frac{M_{ij}}{\lambda}\right)$$



Divergence derived from entropy-regularized cost

• A new divergence

$$D_{\lambda}[p:q] = (1+\lambda) \left(C_{\lambda}(p, \tilde{K}_{\lambda}q) - C_{\lambda}(p, \tilde{K}_{\lambda}p) \right)$$

- Convex with respect to $oldsymbol{p}$ and $oldsymbol{q}$

Proof based on the Riemannian metric $~~{f G}_\lambda=
abla_{m \eta}
abla_{m \eta}arphi_\lambda(m \eta)$



Divergence derived from entropy-regularized cost

• A new divergence

$$D_{\lambda}[p:q] = (1+\lambda) \left(C_{\lambda}(p, \tilde{K}_{\lambda}q) - C_{\lambda}(p, \tilde{K}_{\lambda}p) \right)$$

- Convergence to Wasserstein distance ($\lambda \rightarrow 0$)
- Convergence to a squared distance ($\lambda
 ightarrow \infty$)

$$\lim_{\lambda o \infty} D_{\lambda}[m{p}:m{q}] = rac{1}{2} (m{q}-m{p})^T ilde{m{M}}(m{q}-m{p})$$

 $ilde{m{M}}$: Symmetric matrix determined by $m{M}$

This overcomes disadvantage of C_{λ}

$$\lim_{\lambda \to \infty} C_{\lambda}(\boldsymbol{p}, \boldsymbol{q}) = -\lambda(H(\boldsymbol{p}) + H(\boldsymbol{q}))$$

which cannot measure the distance between $oldsymbol{p}$ and $oldsymbol{q}$

Barycenter of Patterns

 $D_{\lambda}[p:q] = (1+\lambda) \left(C_{\lambda}(p, \tilde{K}_{\lambda}q) - C_{\lambda}(p, \tilde{K}_{\lambda}p) \right)$

• Find the barycenter of input samples p_i (i = 1, 2, ..., N)

$$oldsymbol{q}_C^* = \operatorname*{argmin}_{oldsymbol{q}} \sum_i C_\lambda(oldsymbol{p}_i,oldsymbol{q}) \qquad oldsymbol{q}_D^* = \operatorname*{argmin}_{oldsymbol{q}} \sum_i D_\lambda[oldsymbol{p}_i:oldsymbol{q}]$$

Discrete diffusion operator

Sharpening Conditional metric

$$\boldsymbol{q}_{C}^{*} = \tilde{\boldsymbol{K}}_{\lambda} \boldsymbol{q}_{D}^{*} \qquad \tilde{\boldsymbol{K}}_{\lambda} = \left[\frac{K_{ji,\lambda}}{K_{j\cdot}}\right]_{ij}$$

 $\lambda \to 0$



Under certain conditions,

 $q_D^* = ilde{K}_{\lambda}^{-1} q_C^*$; sharper barycenter by "Inverse diffusion"

Barycenter of Patterns

$$D_{\lambda}[p:q] = (1+\lambda) \left(C_{\lambda}(p, \tilde{K}_{\lambda}q) - C_{\lambda}(p, \tilde{K}_{\lambda}p) \right)$$

Shape-location separation

 $m{q}_D^* = \operatorname*{argmin}_{m{q}} \sum_i D_\lambda[m{p}_i:m{q}]$ is located at the barycenter of $m{p}_i$'s locations





Remark: λ-Divergence [AKO'18]

$$\widetilde{D}_{\lambda}[\boldsymbol{p}:\boldsymbol{q}] = \varphi_{\lambda}(\boldsymbol{p},\boldsymbol{p}) - \varphi_{\lambda}(\boldsymbol{p},\boldsymbol{q}) - \nabla_{\boldsymbol{q}}\varphi_{\lambda}(\boldsymbol{p},\boldsymbol{q}) \cdot (\boldsymbol{p}-\boldsymbol{q})$$

- *Canonical divergence* from the Legendre duality

- It does not satisfy sharpening & shape-location separation

Barycenter of Patterns



 D_{λ} -barycenters are obtained by a gradient method

Summary

- The set of entropy-regularized OT plans is an exponential family and naturally defines the dually flat manifold
 - We obtained its Riemannian metric & dual affine connections
- A new divergence D_{λ}
 - Inherits essential properties of Wasserstein geometry Sharpening, shape-location separation
 - Better than C_{λ}

satisfying criterion of divergence, shaper barycenter

Future work

Information geometry of Entropy-regularized OT with continuous measures

- continuous p vs. continuous q, discrete p vs. continuous q
- Generative model