Sparse Regularized Optimal Transport with Deformed q-Entropy

March 15th, 2023 包含 (Han Bao) / Kyoto University @Workshop OT

つつみ ふくむ Short bio | Han Bao / 包 含

2017 April - 2022 March: Graduate student @ The University of Tokyo

✤ My ex-office was in Science Bldg. #7 (very close from here!)

2022 April -: Project assistant professor @ Kyoto University (Hakubi center)

Research field: computer science, machine learning

Design of robust loss functions

Theory of representation learning











Optimal transport in machine learning

Measure the discrepancy of probability distributions (statistical inference)



• Compute alignments between objects (for information retrieval)









今日紹介する研究

B. & Sakaue. Sparse Regularized Optimal Transport with Deformed q-Entropy. Entropy, 2022.

Summary in one page

• What we do

Propose a new (entropic) regularizer to enhance sparsity of optimal transport solution



What's interesting from practical viewpoint
 Controllable trade-off between sparsity and computational time
 What's interesting from theoretical viewpoint

Connection to convex duality and statistical mechanics via Fenchel–Legendre transform



Introduction Optimal transport

Goal: to find an optimal way to transport goods from sources to targets





satisfy demand constraints



Introduction | Regularized optimal transport

- Motivation: to accelerate optimization (by using convex optimization)
- Primal formulation [Cuturi, 2013]

$$\inf_{\Pi \in U(\mathbf{a},\mathbf{b})} \langle \mathbf{D}, \Pi \rangle + \lambda \sum_{i,j} H(\Pi_{ij}) \quad \text{where } H(\pi) =$$

Sinkhorn algorithm solves the primal by block coordinate descent with $O(n^2)$ time \bigcirc

- Issue: the regularized solution deviate from the original solution (?)
 - In particular: "sparsity" is lost

Cuturi, M. (NeurIPS2013). Sinkhorn distances: Lightspeed computation of optimal transport.

 $\pi \log \pi - \pi$ (negative Shannon entropy)

LP solution 25 10 -15 25

Solution of Sinkhorn







Introduction | Sparsity of transport matrix

- Application of optimal transport: alignment
- Example: matching two different sentence pairs [Swanson et al., 2020]



Swanson, K., Yu, L., & Lei, T. (ACL2020). Rationalizing Text Matching: Learning Sparse Alignments via Optimal Transport.

Sparse alignment is more salient to extract meaningful matches





Motivation | New regularizer to balance sparsity/runtime

```
\langle \mathbf{D}, \mathbf{\Pi} \rangle
         inf
\Pi \in U(\mathbf{a}, \mathbf{b})
```











Solution of Sinkhorn



Quadratic time Dense solution



Formulation | Regularized OT & dual

Regularized OT (primal)

 $\inf_{\mathbf{\Pi}\in U(\mathbf{a},\mathbf{b})} \langle \mathbf{D},\mathbf{\Pi} \rangle + \sum_{i,j} \Omega(\mathbf{\Pi}_{ij}) \text{ for a convex regularizer } \Omega: \mathbb{R} \to \mathbb{R}$

Regularized OT (dual)

$$\sup_{\alpha,\beta\in\mathbb{R}^n} -\langle \mathbf{a},\alpha\rangle - \langle \mathbf{b},\beta\rangle - \sum_{i,j} \Omega^* (-\mathbf{D}_{ij} - \alpha_i - \beta_i)$$

 \Rightarrow By Lagrangian (α, β : multipliers) & strong duality

✤ Note: dual is unconstrained

Primal-dual correspondence (inverse link)

$$\mathbf{\Pi}_{ij} = \nabla \Omega^{\star} \big(- \mathbf{D}_{ij} - \boldsymbol{\alpha}_i - \boldsymbol{\beta}_j \big) \text{ for all i, j}$$

 \Rightarrow Primal transport plan Π can be obtained from dual variables α, β

where Ω^{\star} is Fenchel–Legendre transform $\boldsymbol{\beta}_i)$

Fenchel–Legendre transform

 $\Omega^{\star}(\eta) = \sup \langle \pi, \eta \rangle - \Omega(\pi)$ $\pi \in \operatorname{dom}(\Omega)$





Examples



Regularizer Ω(
$$\pi$$
) = $\frac{\lambda}{2}\pi^2$

♦ Inverse link $\nabla \Omega^{\star}(\eta) = \frac{1}{\lambda} \max\{0, \eta\}$

 Primal-dual correspondence $\mathbf{\Pi}_{ii} = \nabla \Omega^{\star} (-\mathbf{D}_{ii} - \boldsymbol{\alpha}_i - \boldsymbol{\beta}_i)$





Idea | Regularizer from alternative (inverse) link function

ullet Instead of a "dense" inverse link \cdots

$$\mathbf{\Pi}_{ij} = \exp\left(-\frac{\mathbf{D}_{ij} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j}{\lambda}\right)$$

Design a "sparse" inverse link

$$\mathbf{\Pi}_{ij} = \left(-\frac{\mathbf{D}_{ij} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j}{\lambda} \right)$$

Then, recover an entropy-like function

By convex duality (see right figure)

• Q. How to design a sparse inverse link?





Idea q-exponential functions

q-exponential

Used in Tsallis statistics, robust Bayesian inference, etc. Proposed: replace Gibbs kernel with q-exponential distributions

Before:
$$\Pi_{ij} = \exp\left(-\frac{\mathbf{D}_{ij} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j}{\lambda}\right)$$
 (support is
After: $\Pi_{ij} = \exp_q \left(-\frac{\mathbf{D}_{ij} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j}{\lambda}\right)$ (support i

- $\leq q < 1$ = 1

 $\mathbb{R}_{\geq 0}$)

is subset of $\mathbb{R}_{>0}$)





Proposed | Regularizer based on q-exponential functions

Derive primal regularizer from q-exp

$$\Omega(\pi) = \frac{\lambda}{2-q} (\pi \log_q(\pi) - \pi)$$

 \diamond λ : regularizer strength

 \Rightarrow q: control support of kernel

Solve dual optimization

$$\sup_{\boldsymbol{\alpha},\boldsymbol{\beta}\in\mathbb{R}^n} -\langle \mathbf{a},\boldsymbol{\alpha}\rangle - \langle \mathbf{b},\boldsymbol{\beta}\rangle - \sum_{i,j} \Omega^{\star} (-\mathbf{D}_{ij} - \boldsymbol{\alpha}_i - \boldsymbol{\beta}_j)$$

Obtain primal solution via inverse link

$$\mathbf{\Pi}_{ij} = \exp_q \left(-\frac{\mathbf{D}_{ij} + \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j}{\lambda} \right)$$



Convergence analysis

• Suppose dual optimization $\sup_{\alpha,\beta\in\mathbb{R}^n} - \langle \mathbf{a},\alpha\rangle - \langle \mathbf{b},\beta\rangle$

Theorem. Under some conditions, the gradient at the point obtained by *K*-th iteration of BFGS is upper-bounded by

where the rate satisfies 0 < r < 1 and τ is a constant in (0,1).

 \bullet Implication: larger q is beneficial in terms of convergence

$$\langle \boldsymbol{\beta} \rangle - \sum_{i,j} \Omega^{\star} (-\mathbf{D}_{ij} - \boldsymbol{\alpha}_i - \boldsymbol{\beta}_j)$$
 via BFGS

 $\sqrt{\frac{Cn\tau^q}{\lambda}}r^K$





Experiments | Sparsity of transport matrices

• Solution becomes closer to LP solution as $q \rightarrow 0$





Quantitative results

 \diamond q = 1: fully dense (sparsity = 0)

 \diamond q < 1: far more sparse than q = 1 (often sparsity > 0.6)



Sinkhorn (q=1)





Experiments | **Runtime comparison**

• Smaller q often requires more runtime

Trade-off between sparsity and runtime



Dataset size = 100



Dataset size = 300



Comparison of q-entropy and Tsallis entropy

Primal

$$\inf_{\mathbf{\Pi}\in U(\mathbf{a},\mathbf{b})} \langle \mathbf{D},\mathbf{\Pi} \rangle + \sum_{i,j} \Omega(\mathbf{\Pi}_{ij})$$

• q-entropy





 \diamond Inverse link is finitely supported for q < 1

Tsallis entropy

$$T(\pi) = \lambda \pi^q \log_q(\pi)$$

Inverse link

Inverse link is not finitely supported for any q







Summary

• What we do

Propose a new (entropic) regularizer to enhance sparsity of optimal transport solution



Practically: trade-off between sparsity and runtime

Theoretically: design of a good regularizer/entropy/divergence/etc.

Many machine learning problems have convex duality structures!

