Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Optimal transport and statistical inference

Nguyen Xuan Long

University of Michigan

Vietnam Institute for Advanced Studies in Mathematics Hanoi, 7/6/2022

Acknowledgement: Aritra Guha, Jiacheng Zhu, Dat Do, Ding Zhao, Yun Wei, Nhat Ho, Yang Chen, Bach Viet Do, Sunrit Chakraborty

Optimal transport 00000000

Quantifying abstract dependence 00000000

Outline

Data and distances

Optimal transport

Domain adaptation

Distance of latent structures Inverse bounds

Quantifying abstract dependence

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 差 = 釣�?

Optimal transport

Domain adaptation

Quantifying abstract dependence

Data are vast, diverse, complex



- (a) Observations of photon sources near the center of the Orion Nebula from the *Chandra X-ray Observatory* (Jones et al, 2015)
- (b) Trajectories of traffic primitives extracted from sensors-equipped vehicles driven in and around Ann Arbor, Michigan (Guha et al, 2020)
- (c) Latent topics extracted from text corpus

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

From data to statistical models



each dot represents a data point (photon source, car's signal, text doc) we may assume

data ~
$$i.i.d. P(x|\theta)$$

- data are samples of a random variable X
- P is a probability distribution on some observable domain
- unknown θ parameterizes P

Optimal transport

Domain adaptation

Quantifying abstract dependence

Distances on space of distributions

All statistical learning algorithms involve making movements on some (explicit or implicit) space of probability measures

 traditional notions of distance assume existence of density functions (typically wrt Lebesgue measure on Euclidean domain, or counting measure on discrete domain)

Standard distances in statistics, information theory, learning theory

- total variation: $D_{TV}(p,q) = \frac{1}{2} \int |p(x) q(x)| d\mu$
- Hellinger:

$$h(p,q) = \frac{1}{\sqrt{2}} \|\sqrt{p} - \sqrt{q}\|_{L^2}$$

• relative entropy (KL divergence):

$$\mathcal{K}(p,q) = \int p(x) \log(p(x)/q(x)) d\mu$$

we also know

$$h^2 \leq D_{TV} \leq \sqrt{2}h \leq \sqrt{K}.$$

Optimal transport •0000000 Domain adaptation

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Outline

Data and distances

Optimal transport

Domain adaptation

Distance of latent structures Inverse bounds

Quantifying abstract dependence

Optimal transport

Domain adaptation

Quantifying abstract dependence 00000000

There may be some metric structure in the supports of P and Q, but they may be disjoint



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

Optimal transport 0000000



Monge problem Find a map $x \mapsto T(x)$ s.t. if $X \sim P$ then $Y := T(X) \sim Q$. That is, find T such that the pushforward measure satisfies $T_{\#}P = Q$.

Kantorovich problem enlarges Monge's into a solvable problem: find a "stochastic map", i.e., a coupling of P and Q that minimizes expected cost of transportation ▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Optimal transport

Domain adaptation

Quantifying abstract dependence 00000000

Basic definition and facts

Let (\mathcal{X}, μ) and (\mathcal{Y}, ν) be two complete and separable metric probability spaces. Let $c : \mathcal{X} \times \mathcal{Y} \to \mathbb{R} \cup \{+\infty\}$ be a non-negative lower-semicontinuous cost function.

<u>Lemma</u> Let $\Pi(\mu, \nu)$ be the space of all <u>couplings</u> of μ, ν , i.e., all joint distributions on $\mathcal{X} \times \mathcal{Y}$ that admit marginal distributions μ and ν . Then there exists a coupling $\mu \in \Pi(\mu, \nu)$ that minimizes the total cost

$$\mathbb{E}_{\pi}c(X,Y)=\int c(x,y)d\pi(x,y).$$

Optimal transport

•

Domain adaptation

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Kantorovich duality

$$\begin{split} \min_{\pi \in \Pi(\mu,\nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x,y) d\pi(x,y) = \\ \max_{(\psi,\phi) \in L^1(\mu) \times L^1(\nu): \phi + \psi \le c} \int \psi(x) d\mu(x) + \phi(y) d\nu(y). \end{split}$$

- (i) \geq is trivial; = is due to convex optimization in Banach spaces
- (ii) dual formulation allows precise characterization of the optimal π and corresponding optimal φ, ψ: assume that the optimal cost is finite, then the support of π satisfies ψ(x) + φ(y) = c(x, y) almost surely.

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

Kantorovich duality

$$\begin{split} \min_{\pi\in\Pi(\mu,\nu)} \int_{\mathcal{X}\times\mathcal{Y}} c(x,y) d\pi(x,y) = \\ \max_{(\psi,\phi)\in L^1(\mu)\times L^1(\nu):\phi+\psi\leq c} \int \psi(x) d\mu(x) + \phi(y) d\nu(y). \end{split}$$

- (i) \geq is trivial; = is due to convex optimization in Banach spaces
- (ii) dual formulation allows precise characterization of the optimal π and corresponding optimal φ, ψ: assume that the optimal cost is finite, then the support of π satisfies ψ(x) + φ(y) = c(x, y) almost surely.
- (iii) when c is continuous, then the support of π is c-cyclic monotone, i.e., the set of pairs (x_n, y_n) such that

$$\sum_{i=1}^{N} c(x_i, y_i) \leq \sum_{i=1}^{N} c(x_i, y_{i+1})$$

hold for all such pairs and all N (with the convention $y_{N+1} = y_1$).

・ロト・西ト・山田・山田・山口・

Optimal transport

Domain adaptation

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Two useful properties

(1) Wasserstein metric: Let (X, d) be a Polish metric space, $r \in [1, \infty)$. For any two pm's μ, ν on \mathcal{X} , the Wasserstein metric of order r is given given by

$$W_r(\mu,\nu) = \left(\inf_{\pi\in\Pi(\mu,\nu)}\int_{\mathcal{X}} d(x,y)^r d\pi(x,y)\right)^{1/r}$$

= $\inf\left\{ [\mathbb{E}d(X,Y)^r]^{1/r}, X \sim \mu, Y \sim \nu \right\}.$

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

Two useful properties

(I) Wasserstein metric: Let (X, d) be a Polish metric space, $r \in [1, \infty)$. For any two pm's μ, ν on \mathcal{X} , the Wasserstein metric of order r is given given by

$$W_r(\mu,\nu) = \left(\inf_{\pi\in\Pi(\mu,\nu)}\int_{\mathcal{X}} d(x,y)^r d\pi(x,y)\right)^{1/r}$$

= $\inf\left\{ [\mathbb{E}d(X,Y)^r]^{1/r}, X \sim \mu, Y \sim \nu \right\}.$

We can define properly distance metric on distribution of any "complex" data populations, as long as we have a metric d on the data instances available

Data and distances	Optimal transport	Domain adaptation	Distance of latent structures	Quantifying abstract dependence
0000	00000000	000000000000000000000000000000000000000	0000000 0000000000000000000000000000000	0000000

(II) Brenier-Rachev-Ruschendorf theorem: Let $c(x, y) = ||x - y||^2$ in \mathbb{R}^n . μ and ν two pm's with bounded second moments. If μ is absolute continuous wrt the Lebesgue measure, then there is a unique optimal coupling of π of x, y, under which y is <u>uniquely</u> determined by x almost surely. In fact, for some lsc convex function ψ ,

 $y \in \partial \psi(x)$ almost surely.

In other words, the Monge transport plan $x \mapsto T(x) = \partial \psi(x)$ exists, and the optimal coupling π is characterized as the pushforward measure from the source distribution μ :

$$\nu = T_{\#}\mu$$

$$\pi = (Id, T)_{\#}\mu.$$

See Villani (2008) and Ambrosio-Gigli-Savare (2005) for much recent advances in analysis, PDEs, differential geometry associated with OT

Domain adaptation

Quantifying abstract dependence

Implications to statistical learning/ inference

- (1) we can define distance between any data populations, which inherit the metric d of data instances
- (2) if we know the source distribution for X ~ μ, we can perhaps model the target Y ~ ν via a (optimal transport) map X → Y = T(X), i.e., ν is a pushforward measure of μ by the map T:

$$\nu = T_{\#}(\mu)$$

- the question of "learning" a target distribution becomes the learning of transport map T, e.g., Wasserstein-GAN (Arjovsky et al, 2017)
- (3) while all this seems nice, optimal transport is only useful if the metric on the data d is meaningful, and map T can be effectively learned
 - we will illustrate this in the following "domain adaptation problem"

Optimal transport 00000000 Domain adaptation
OOOOOOOOOOOOO

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Outline

Data and distances

Optimal transport

Domain adaptation

Distance of latent structures Inverse bounds

Quantifying abstract dependence

Quantifying abstract dependence 00000000

Domain adaptation problems

1. If I learned to drive in Arizona, can I adapt my experience to driving in California? Answer: yes.

How about driving in Hanoi?

- here, driving experience = samples of time series of driving trajectories
- 2. If a robot knows how to pick up an object, how can it be taught to bend and pick up efficiently?

General problem If we know/ can learn well distribution μ , how can we adapt μ to learn another distribution ν

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

・ロト ・ 国 ト ・ ヨ ト ・ ヨ ト

Our approach: functional optimal transport

Main ideas: (Zhu, Guha, Do, Xu, Nguyen and Zhao, 2021)

- each data point is a realization of a random function
- given source and target distributions μ and ν on functions
- learn the transport map from space of compact linear operators

Pushing forward sampled paths for source to target distribution



Geodesic linking source distribution to target distribution



Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

Learning OT map in function space

Let H_1, H_2 be Hilbert space of functions endowed with Borel probability measures μ_1 and μ_2 , resp.

Let $\mathcal{B}_{HS}(H_1, H_2)$ be the space of Hilbert-Schmidt operators, i.e., a Hilbert space of linear operators endowed with the scalar product

$$\langle A, B \rangle_{HS} = \sum_{i=1}^{\infty} \langle AU_i, BU_i \rangle_{H_2}$$

where $(U_i)_{i=1}^{\infty}$'s form a complete orthonormal basis of H_1 .

Consider the optimization problem

$$\inf_{T \in \mathscr{B}_{HS}(H_1, H_2)} J(T) := W_2^2(T_{\#}\mu, \nu) + \eta \|T\|_{HS}^2$$
(1)

 $\eta > 0$ is a regularization parameter.

Optimal transport

Domain adaptation

Quantifying abstract dependence 00000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

Lemma

Assume that μ and ν have bounded second moments:

$$E_{f_1 \sim \mu} \|f_1\|_{H_1}^2 < \infty, \quad E_{f_2 \sim \nu} \|f_2\|_{H_2}^2 < \infty$$
 (A.1)

then the objective (1) is a strictly convex function, which admits a unique minimizer.

Optimal transport 00000000 Domain adaptation 0000000000000

Quantifying abstract dependence 00000000

In practice, given i.i.d. samples

$$f_{11}, f_{12}, \dots, f_{1n_1} \sim \mu,$$

 $f_{21}, f_{22}, \dots, f_{2n_2} \sim \nu,$

the empirical version of our optimization problem becomes:

$$\inf_{T \in \mathscr{B}_{HS}} \hat{J}_n(T), \quad \hat{J}_n(T) := W_2^2(T_{\#}\hat{\mu}_{n_1}, \hat{\nu}_{n_2}) + \eta \|T\|_{HS}^2, \tag{2}$$

where $\hat{\mu}_{n_1} = \frac{1}{n_1} \sum_{l=1}^{n_1} \delta_{f_{1l}}$ and $\hat{\nu}_{n_2} = \frac{1}{n_2} \sum_{k=1}^{n_2} \delta_{f_{2k}}$ are the empirical measures, and $n = (n_1, n_2)$.

Moreover, restrict T to a $K_1 \times K_2$ dimensional subspace of \mathscr{B}_{HS}

うしん 前 ふかく ボット 間 うくの

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

Lemma Under assumption (A.1), the following hold.

For any fixed C₀ > 0,

$$\sup_{|T|| \le C_0} |\hat{J}_n(T) - J(T)| \xrightarrow{P} 0 \quad (n \to \infty).$$
(3)

For any n, K, Ĵ_n has a unique minimizer T̂_{K,n} over B_K. Moreover, there exists a finite constant D such that P(sup_K || T̂_{K,n} || < D) → 1 as n → ∞.

It then follows that

Theorem

The minimizer of Eq. (2) for $\hat{T}_{K,n} \in B_K$ is a consistent estimate for the minimizer of Eq. (1). Specifically, $\hat{T}_{K,n} \xrightarrow{P} T_0$ as $K_1, K_2, n_1, n_2 \to \infty$.

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Additional considerations

- (i) in practice, sampled functions $(f_{1l})_{l=1}^{n_1}$, $(f_{2l})_{l=1}^{n_2}$ are observed only at a finite number of design points d_1, d_2 , resp.
- (ii) the HS operator is approximated by dimension truncation (to spaces of $K_1 \times K_2$ matrices)
- (iii) consistency theory can be extended under the regime that $n_1, n_2, K_1, K_2 \rightarrow \infty$, and $d \rightarrow \infty$ suitably
 - see the paper (Zhu et al, 2021) for details.

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

FOT vs non-functional approaches



Optimal transport

Domain adaptation

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

FOT vs non-functional approaches: quantitative comparison

Method	1 ightarrow 1	1→2	2→1	2→2	2→3
GPOT	17.560	12.895	15.263	61.561	39.159
LSOT	133.434	94.229	117.832	929.108	663.461
DSOT	6.871	13.226	9.679	46.521	41.009
FOT	2.873	11.982	3.316	44.071	32.547

Table: Quantitative comparison on the mixture of sinusoidal functions data. The maps obtained by FOT method achieved the best performance under the Wasserstein distance objective.

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Domain adaptation for robotic arm movements



(a) The arm of the Baxter robot and the Sawyer robot used in MIME dataset and Roboturk dataset. They share a similar structure, 7 joints and one end effector.



(b) Source motion: "Roboturk-bins-Bread" by Sawyer.



(c) Target motion: "MIME Picking (left-hand)" by Baxter.



(d) The pushforward motion of the transport map looks like Baxter's but inherits trait of Sawyer's motions.

Optimal transport

Domain adaptation 00000000000000

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Augmenting deep learning models

Semi-supervised learning

Time series motion prediction tasks with FOT augmentation (use FOT to generate more samples in target domain for training with deep learning predictive models)

Method	LSTM	ANP	RANP	MAML*	TL*	FOTLSTM	FOT _{ANP}	FOT _{RANP}	FOT _{MAML}	FOT _{TL}
R1→M1	2.0217	1.3261	1.9874	0.0307	0.5743	0.0271	0.0963	0.0687	0.0165	0.0277
$R1 \rightarrow M2$	1.6821	1.0951	1.5681	0.0374	0.7083	0.0414	0.1642	0.1331	0.0191	0.0446
$R2 \rightarrow M1$	1.3963	0.6642	1.7256	0.0327	0.2491	0.0277	0.0951	0.0696	0.0202	0.0906
$R2 \rightarrow M2$	1.1952	0.6307	1.3659	0.0477	0.4020	0.0331	0.1620	0.1554	0.0167	0.0406

Table: MSE error results of different predictive models.

R1: Roboturk-bins-bread, R2: Roboturk-pegs-RoundNut, M1:MIME1-Pour-left, M2: MIME12-Picking-left.

I distances Optimal tran 0000000

timal transport

Domain adaptation

Quantifying abstract dependence 00000000

Our story so far: using the pushforward map T to model a target distribution via

$$u = T_{\#}\mu$$

in a complex data domains, but there are huge technical challenges, when

- incorporating domain knowledge of the support of ν , and T,
- and when the interest is *not* data distribution but on distribution on meaningful quantities related to it

This leads to optimal transport on distributions on the space of quantities of interest. Thus, the starting point is with probability models via latent *random* variables/structures.

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

▲ロ▶ ▲周▶ ▲ヨ▶ ▲ヨ▶ ヨ のなべ

Outline

Data and distances

Optimal transport

Domain adaptation

Distance of latent structures Inverse bounds

Quantifying abstract dependence

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence 00000000

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

Let's talk about latent structured models!

Optimal transport

Domain adaptation 0000000000000000 Quantifying abstract dependence 00000000

Observations of the Orion constellation

ancient data



an ancient model



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Optimal transport

Domain adaptation

Quantifying abstract dependence 00000000

Observations of the Orion constellation

ancient data



Orion nebula



an ancient model



X-ray data via Chandra observatory



Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

Consider a situation where the quantity of interest is not data distribution itself: photon sources in the Orion Nebula



Assume the pixel locations are sample from a mixture distribution

$$X_1,\ldots,X_n\sim\sum_{j=1}^K p_j f(x|\theta_j)$$

there are K photon sources; θ_j represents information about the arrival time, location, energy level of photon source j

- *p_j*: the probability that the photon comes from source *j* for certain type of stellar and inter-galatic events of interest (exploded stars, star birth, etc)
- f is a probability kernel which captures distribution of observations (e.g., King profile for the location observations)

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

The mixture distribution used to describe the observed locations has very little scientific meaning:

$$\mathcal{P}_{\mathcal{G}} = \sum_{j=1}^{K} p_j f(x| heta_j)$$

Of interest to astrophysicists and astronomers is the mixing measure

$${\it G} = \sum_{j=1}^{K} {\it p}_j \delta_{ heta_j}$$

Here, optimal transport continues to play a fundamental role in characterizing the learning behavior of quantities of interest, where the underlying metric structure of support is derived from the rich structure of the probabilistic models

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ● ●

Wasserstein metric on G

How to define a metric on the space of θ 's:

- suppose $heta \in \mathbb{R}^d$, we may take $d(heta, heta') := \| heta heta'\|_r^r$
- better yet, make use of Hellinger distance:

$$d(\theta, \theta') := h(f(\cdot|\theta), f(\cdot|\theta'))$$

• this results in the following Wasserstein metric

$$W(G,G') = \inf_{\pi \in \Pi(G,G')} \int d(\theta,\theta') d\pi$$

This is called "composite distance" in Nguyen (2013)

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Statistical learning methods

<u>Statistical formulation</u>: given $X_1, \ldots, X_n \stackrel{\text{i.i.d.}}{\sim} P_G$ for some "truth"

$$G = G_0 = \sum_{j=1}^{K_0} p_j^0 \delta_{\theta_j^0}$$

To learn G, we can either apply

• maximum likelihood estimate via the EM algorithm:

$$\hat{G} := \operatorname{argmax}_G \sum_{i=1}^n \log p_G(X_i)$$

• Bayesian method: place a prior distribution on G, and apply Bayes formula to obtain the posterior dist. $\Pi(G|X_1, \ldots, X_n)$.

Theoretical questions: in what sense does the estimate \hat{G} , or the posterior distribution for G converge to the truth G_0 in the Wasserstein space?
Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Statistical learning methods

<u>Statistical formulation</u>: given $X_1, \ldots, X_n \stackrel{\text{i.i.d.}}{\sim} P_G$ for some "truth"

$$G = G_0 = \sum_{j=1}^{K_0} p_j^0 \delta_{\theta_j^0}$$

To learn G, we can either apply

maximum likelihood estimate via the EM algorithm:

$$\hat{G} := \operatorname{argmax}_G \sum_{i=1}^n \log p_G(X_i)$$

• Bayesian method: place a prior distribution on G, and apply Bayes formula to obtain the posterior dist. $\Pi(G|X_1, \ldots, X_n)$.

Theoretical questions: in what sense does the estimate \hat{G} , or the posterior distribution for G converge to the truth G_0 in the Wasserstein space?

A key to these questions is the derivation of inverse bounds a = b = a

Optimal transport

Domain adaptation

 Quantifying abstract dependence

A simpler problem: deconvolution problem

Example: we are interested in the distribution of signal $\theta_i \in \mathbb{R}^d$, i = 1, ..., n, given i.i.d. noisy observations of the form

$$\begin{array}{lll} X_i &=& \theta_i + \epsilon_i, \\ \epsilon_i &\sim \mathrm{i.i.d.} & f, \ \epsilon_i \perp \theta_i \end{array}$$

Suppose that $\theta_i \sim \text{i.i.d.} G_0$, then $X_i \sim \text{i.i.d.} f * G_0$.

Optimal transport

Domain adaptation

Distance of latent structures

Quantifying abstract dependence

A simpler problem: deconvolution problem

Example: we are interested in the distribution of signal $\theta_i \in \mathbb{R}^d$, i = 1, ..., n, given i.i.d. noisy observations of the form

$$\begin{array}{lll} X_i &=& \theta_i + \epsilon_i, \\ \epsilon_i &\sim \mathrm{i.i.d.} & f, \ \epsilon_i \perp \theta_i \end{array}$$

Suppose that $\theta_i \sim \text{i.i.d.} G_0$, then $X_i \sim \text{i.i.d.} f * G_0$.

An inverse bound is an inequality of the type

$$W_2(G_0,G) \leq \Phi(V(f * G_0, f * G))$$

where

- LHS is an optimal transport distance e.g., W_1, W_2, \ldots on distribution of the latent θ
- V is the total variation distance of data populations $f * G_0$ and f * G
- $\Phi: [0,\infty) \to [0,\infty)$ is a strictly increasing function, $\Phi(0) = 0$.

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ・ つ へ ()・

Optimal transport 00000000 Domain adaptation 00000000000000 Distance of latent structures

Quantifying abstract dependence

why inverse bounds are important

- Inverse bounds are useful for deriving rates of convergence of the latent θ distributions, once the convergence in distribution of data population X has been established
- The opposite direction of the inverse bound is typically easy
 - let $(heta, heta_0)$ be an optimal coupling of G and G_0
 - then (X, X_0) where $X = \theta + \epsilon$, $X_0 = \theta_0 + \epsilon$ represent a coupling of f * G and $f * G_0$, so

$$W_2^2(f * G, f * G_0) \leq \mathbb{E} \|X - X_0\|^2 = W_2^2(G, G_0).$$

• to obtain bound for $V(f * G, f * G_0)$, by an application of Jensen's inequality one gets

$$\begin{aligned} \|f * G - f * G_0\|_{L^1} &\leq \quad \mathbb{E}_{\theta,\theta_0} \|f(\cdot - \theta) - f(\cdot - \theta_0)\|_{L^1} \\ &\leq \quad CW_1(G,G_0), \end{aligned}$$

where C is an "integrated" Lipschitz constant of the density function $f(\cdot - \theta)$ wrt θ .

 similar upper bounds on any f-divergence using the same technique (Nguyen, 2013)

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Inverse bounds for convolutional models $X = \theta + \epsilon$

- let f be a pdf on \mathbb{R}^d that is symmetric around 0, and the Fourier transform of f satisfies $\tilde{f}(\omega) \neq 0$ for all $\omega \in \mathbb{R}^d$.
- Borrowing from deconvolution literature (cf. Fan (1991)) say
 - f is ordinary smooth with parameter $\beta > 0$ if $\int_{[-1/\delta, 1/\delta]^d} \tilde{f}^{-2} d\omega \lesssim (1/\delta)^{2d\beta}$ as $\delta \to 0$,
 - f is supersmooth with parameter $\beta > 0$ if $\int_{[-1/\delta, 1/\delta]^d} \tilde{f}^{-2} d\omega \lesssim \exp(2d\delta^{-\beta})$ as $\delta \to 0$.

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Inverse bounds for convolutional models $X = \theta + \epsilon$

- let f be a pdf on \mathbb{R}^d that is symmetric around 0, and the Fourier transform of f satisfies $\tilde{f}(\omega) \neq 0$ for all $\omega \in \mathbb{R}^d$.
- Borrowing from deconvolution literature (cf. Fan (1991)) say
 - f is ordinary smooth with parameter $\beta > 0$ if $\int_{[-1/\delta, 1/\delta]^d} \tilde{f}^{-2} d\omega \lesssim (1/\delta)^{2d\beta}$ as $\delta \to 0$,
 - f is supersmooth with parameter $\beta > 0$ if $\int_{[-1/\delta, 1/\delta]^d} \tilde{f}^{-2} d\omega \lesssim \exp(2d\delta^{-\beta})$ as $\delta \to 0$.
- for any pair of G and G' whose support lie in a bounded subset of \mathbb{R}^d , then (Nguyen, 2013)
 - if f is ordinary smooth, then for any $m < 4/(4 + (2\beta + 1)d)$, for some constant $C(d, \beta, m)$,

$$W_2^2(G,G') \leq C(d,\beta,m)V(f*G,f*G')^m.$$

• if f is supersmooth, then there is $C(d,\beta) > 0$ such that

$$W_2^2(G,G') \leq C(d,eta) [-\log V(f*G,f*G')]^{-2/eta}.$$

うとの 川田 (中国)・ (田)・ (日)・

Optimal transport 00000000 Domain adaptation 0000000000000 Distance of latent structures

Quantifying abstract dependence

Beyond convolutional model

- let kernel f be the Gaussian pdf on \mathbb{R} : $f(x|\mu, \nu)$.
- let G be a distribution on the bivariate parameter $(\mu, \nu) \in \mathbb{R} \times \mathbb{R}_+$.
- then the marginal distribution of the observed X is the mixture distribution

$$X \sim P_G := \int f(\cdot|\mu,\nu) dG(\mu,\nu)$$

• Open question: obtain an inverse bound of the form

$$W_2(G_0,G) \leq \Phi(V(P_{G_0},P_G))$$

where the optimal transport distance is defined in a natural way for the distribution on $\mathbb{R}\times\mathbb{R}_+$, under the metric, e.g.:

$$d((\mu_1, \nu_1), (\mu_2, \nu_2)) := |\mu_1 - \mu_2| + |\nu_1^{-1} - \nu_2^{-1}|.$$

Optimal transport 00000000 Domain adaptation 00000000000000 Distance of latent structures

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Where are the "great" models in our time?



Gato: a "multi-modal, multi-task, multi-embodiment generalist agent"

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

Where are the "great" models in our time?



Gato: a "multi-modal, multi-task, multi-embodiment generalist agent" "...Gato can sense and act with different embodiments across a wide range of environments using a single neural network with the same set of weights. (It) was trained on 604 distinct tasks..."

Deepmind's publication "A Generalist Agent", 5/19/2022

Distance of latent structures

Nando de Freitas: "The Game is Over!" But is it?





a living room with three different. Aman in a bise suit with a white. Man holding up a banana to color deposits on the floor

a morn with a long red sug a tyand some pictures.



wearing a suit and tie.

box tin and block shoes. A man with a hat in his hand looking at the camera.



A bearded man is holding a plate of food.

take a picture of it.



a man smilles while holding up a Two horses are laying on their slice of cake



a group of people that is next to Man biting a kite while standing a big horse

A tan horse holding a piece of cloth lying on the ground

side on the dirt.



on a construction site

a big truck in the middle of a

A truck with a kite painted on the back is parked by rocks.



a white horse with a blue and silver bridle

A white borne with blue and gold chains.

a wal. wave.



needs

A surfer riding a wave in the **b**dean

A horse is being shown behind A surfer with a wet suit riding a



a couple of people are out in the A baseball player pitching a ball. Patachius on top of a bowl with A group of phildren eating pizza on top of a baseball field

nitcher on a baseball finit

A baseball player at bat and a catcher in the dirt during a baseball game



coffee on the side.

Amen throwing a baseball at a A bowl and a glass of liquid sits. Two bows having pizza for lunch on a table.

> A while plate filled with a banana bread next to a cup of coffee



at a toble.

with their friends.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

The boy's are eating pizza logether at the table

We heard similar claims before: including but not restricted to neural nets circa 1950s, 1980s (convolutional neural nets), 2010s (cnn on steroid), 2020s (now with transformers)

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Despite truly impressive demos, AI models are still <u>unreliable</u> in domains where reliability really matters!

These models don't understand their domain, and neither do we understand if and when they work! But we always try to make progress at creating and (hopefully) understanding them.

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Despite truly impressive demos, AI models are still <u>unreliable</u> in domains where reliability really matters!

These models don't understand their domain, and neither do we understand if and when they work! But we always try to make progress at creating and (hopefully) understanding them.

"great models are temporary, mathematical principles are forever"

more to the point, I mean the mathematics which helps to justify the presence of latent (hidden) variables, the mechanism for memory and attention, and the statistical/computational theory which helps to explain and achieve the emergence of such representations

Domain adaptation

Distance of latent structures

Quantifying abstract dependence

From data assumptions to models of populations

de Finetti's Theorem: If $(X_j)_{j=1}^{\infty}$ is an infinite exchangeable sequence of random variables, i.e.,

$$(X_1,\ldots,X_N) \stackrel{d}{=} (X_{\pi(1)},\ldots,X_{\pi(N)}) \ \forall N, \forall \pi$$

then there exists a random probability P such that

$$X_1, X_2, \ldots | P \stackrel{i.i.d.}{\sim} P$$

Mixtures of product distributions for N-sequence X₁,..., X_N: conditionally given some θ, the X_i are i.i.d.

$$P_{G,N}(X_1 \in A_1,\ldots,X_N \in A_N) = \int \prod_{n=1}^N P_{\theta}(X_n \in A_n|\theta)G(d\theta)$$

- kernel $P_{ heta}$ known and uniquely parameterized by $heta\in\Theta$
- G is de Finetti mixing measure on Θ;
 G characterizes heterogeneity of underlying data population =

Optimal transport

Domain adaptation 00000000000000 Distance of latent structures

Quantifying abstract dependence

Compositional data structures

Data are often composed of a collection of dependent populations

- there are multiple hospitals, each hospital has many patients
- there are different animals, each animal carry a set of genes
- different countries, each of which is organized into regions, each of which is organized into counties, with residents in each of them
- "activity recognition problem": a collection of computer users, each user is associated with a collection of computer related activities (organized by days), each day has a collection of activities (apps run)
- a collection of text corpora, each text corpus is a collection of documents, each document is a collection of words
- a database of images divided by groups, each image is a collection of image patches, each patch a collection of pixels or other specific computer vision elements

Optimal transport 00000000 Domain adaptation 00000000000000 Distance of latent structures

Quantifying abstract dependence 00000000

Exchangeable collection of data sets

Each data set modeled via a mixture model

 \Longrightarrow they are coupled to enable "borrowing of strength"



[courtesy M. Jordan's slides]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

This gives rise naturally to a hierarchical model

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

イロト 不得 トイヨト イヨト

-

A hierarchical model setting

m groups of data, each of which is given an n-sample



A key point here: $\mathscr{D}_{\alpha G}$ represents a distribution on the space of distributions

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Metrics on Bayesian hierarchies of distributions

Need a notion of distance between, say $\mathscr{D}_{\alpha {\it G}}$ and $\mathscr{D}_{\alpha ' {\it G}'}$

Recall: for $G, G' \in \mathcal{P}(\Theta)$, space of Borel probability measures on Θ ,

$$W_r(G,G') := \inf_{\kappa \in \mathcal{T}(G,G')} \left[\int \|\theta - \theta'\|^r d\kappa(\theta,\theta')
ight]^{1/r}$$

 $\mathcal{T}(G,G')$ is the space of all couplings of G,G'.

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

Metrics on Bayesian hierarchies of distributions

Need a notion of distance between, say $\mathscr{D}_{\alpha {\it G}}$ and $\mathscr{D}_{\alpha ' {\it G}'}$

Recall: for $G, G' \in \mathcal{P}(\Theta)$, space of Borel probability measures on Θ ,

$$W_r(G,G') := \inf_{\kappa \in \mathcal{T}(G,G')} \left[\int \|\theta - \theta'\|^r d\kappa(\theta,\theta') \right]^{1/r}.$$

 $\mathcal{T}(G,G')$ is the space of all couplings of G,G'.

Distance between measures of measures in Bayesian hierarchy: Let $\mathcal{D}, \mathcal{D}' \in \mathcal{P}(\mathcal{P}(\Theta))$ (the space of Borel probability measures on $\mathcal{P}(\Theta)$). Define Wasserstein distance between $\mathcal{D}, \mathcal{D}'$

$$W_r(\mathcal{D},\mathcal{D}'):=\inf_{\mathcal{K}\in\mathcal{T}(\mathcal{D},\mathcal{D}')}\left[\int W_r^r(G,G')\ d\mathcal{K}(G,G')\right]^{1/r}$$

 $\mathcal{T}(\mathcal{D},\mathcal{D}')$ is the space of all couplings of $\mathcal{D},\mathcal{D}'\in\mathcal{P}(\mathcal{P}(\Theta))$

・ロト・日本・日本・日本・日本・日本

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence

A lemma and a "theorem" (Nguyen, Bernoulli 2016)

Lemma

(a) Let $G, G' \in \mathcal{P}(\Theta)$, and $\mathcal{D}, \mathcal{D}' \in \mathcal{P}(\mathcal{P}(\Theta))$ such that $\int Pd\mathcal{D} = G$ and $\int Pd\mathcal{D}' = G'$. For $r \geq 1$, if $W_r(\mathcal{D}, \mathcal{D}')$ is finite then

 $W_r(\mathcal{D}, \mathcal{D}') \geq W_r(G, G').$

(b) If
$$\mathcal{D} = \mathscr{D}_{\alpha G}$$
 and $\mathcal{D}' = \mathscr{D}_{\alpha G'}$ (same α), then
 $W_r(\mathcal{D}, \mathcal{D}') = W_r(G, G').$

"Theorem"

- when the number of groups *m* increases and the sample size *n* increases suitably, the posterior distribution of *G* contracts to true G₀ under the above Wasserstein metric
- individual group admits improved posterior contraction due to the "borrowing of information" from other groups

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

Setting and assumptions



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへ⊙

Optimal transport

Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ → □ ・ つくぐ

Setting and assumptions



Setting

- Let Θ be a bounded subset of \mathbb{R}^d .
- True Dirichlet base measure $G_0 \in \mathcal{P}(\Theta)$ is atomic.
- Given *m* Dirichlet processes Q_1, \ldots, Q_m drawn from $\mathcal{D}_{\alpha G}$, for $G = G_0$.
- For each process Q_i there is an *n*-sample from the mixture dist Q_i * f

Optimal transport

Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Assumptions

- (A1) For some $r \ge 1$, $C_1 > 0$, $h(f(\cdot|\theta), f(\cdot|\theta')) \le C_1 \|\theta \theta'\|^r$ and $K(f(\cdot|\theta), f(\cdot|\theta')) \le C_1 \|\theta \theta'\|^r \ \forall \theta, \theta' \in \Theta$.
- (A2) There holds $M = \sup_{\theta, \theta' \in \Theta} \chi(f(\cdot|\theta), f(\cdot|\theta')) < \infty$.
- (A3) *G* is endowed with Dirichlet prior $\mathscr{D}_{\gamma H}$, where $H \in \mathcal{P}(\Theta)$ is non-atomic.

Distance of latent structures

Quantifying abstract dependence

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Posterior concentration theorem (Nguyen, 2016)

As $n \to \infty$ and $m \to \infty$, the posterior distribution of Dirichlet base measure G concentrates to G_0 at the rate

$$\epsilon_{m,n} \asymp \left(\frac{n^{3d}\log m}{m}\right)^{1/(2d+2)} + A(\delta_n)$$

Distance of latent structures

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Posterior concentration theorem (Nguyen, 2016)

As $n \to \infty$ and $m \to \infty$, the posterior distribution of Dirichlet base measure G concentrates to G_0 at the rate

$$\epsilon_{m,n} \asymp \left(\frac{n^{3d}\log m}{m}\right)^{1/(2d+2)} + A(\delta_n).$$

 $\delta_n \rightarrow 0$ is the demixing rate — the rate of estimating mixing measure Q from an *n*-sample of a mixture density Q * f (obtaining this rate is the earlier focus)

Distance of latent structures

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Posterior concentration theorem (Nguyen, 2016)

As $n \to \infty$ and $m \to \infty$, the posterior distribution of Dirichlet base measure G concentrates to G_0 at the rate

$$\epsilon_{m,n} \asymp \left(\frac{n^{3d}\log m}{m}\right)^{1/(2d+2)} + A(\delta_n).$$

 $\delta_n \rightarrow 0$ is the demixing rate — the rate of estimating mixing measure Q from an *n*-sample of a mixture density Q * f (obtaining this rate is the earlier focus)

function A depends on the geometric structure of the support of G_0

Data and distances Optimal transport Domain adaptation Distance of latent structures Quantifying abstract de	rependence
0000 00000000 00000000 0000000 0000000 0000	

(ii) If f is supersmooth with parameter β , then it suffices to set

$$\frac{m}{\log m(\log n)^{\alpha^*(2d+2)/\beta}} \lesssim n^{3d} \lesssim \frac{m}{\log m}$$

In particular, if *n* satisfies $n^{3d} (\log n)^{\alpha^*(2d+2)/\beta} \simeq \frac{m}{\log m}$, then we obtain the concentration rate $\epsilon_{m,n} \simeq (\log n)^{-\alpha^*/\beta} \simeq (\log m)^{-\alpha^*/\beta}$.

(iii) Requirements of the type $n_1(m) \le n \le n_2(m)$ appear crucial in deriving posterior concentration rates in hierarchical models. It is an interesting open question to establish the concentration behavior (or the lack thereof) for the full range of n.

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

(i) if G_0 has a finite and unknown number of support points on a bounded subset of \mathbb{R}^d , then

$$A(\delta_n) \asymp \delta_n^{\alpha^*/(\alpha^*+1)}.$$

where $\alpha^* = \inf_{\theta \in \text{spt } G_0} \alpha G_0(\{\theta\})$,

Optimal transport

Domain adaptation

Distance of latent structures

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

 (i) if G₀ has a finite and unknown number of support points on a bounded subset of R^d, then

$$A(\delta_n) \asymp \delta_n^{\alpha^*/(\alpha^*+1)}.$$

where $\alpha^* = \inf_{\theta \in \text{spt } G_0} \alpha G_0(\{\theta\})$,

(ii) if G_0 has infinite and supersparse support on \mathbb{R}^d ,

$$A(\delta_n) \asymp \exp - [\log(1/\delta_n)]^{1/(1 \lor \gamma_0 + \gamma_1)}$$

(iii) if G_0 has infinite and ordinary sparse support on \mathbb{R}^d ,

$$A(\delta_n) \asymp [\log(1/\delta_n)]^{-1/(\gamma_0+\gamma_1)}$$

Optimal transport 00000000 Domain adaptation

Distance of latent structures

Quantifying abstract dependence 00000000

(i) if G_0 has a finite and unknown number of support points on a bounded subset of \mathbb{R}^d , then

$$A(\delta_n) \asymp \delta_n^{\alpha^*/(\alpha^*+1)}.$$

where $\alpha^* = \inf_{\theta \in \text{spt } G_0} \alpha G_0(\{\theta\})$,

(ii) if G_0 has infinite and supersparse support on \mathbb{R}^d ,

$$A(\delta_n) \asymp \exp - [\log(1/\delta_n)]^{1/(1 \lor \gamma_0 + \gamma_1)}$$

(iii) if G_0 has infinite and ordinary sparse support on \mathbb{R}^d ,

$$A(\delta_n) \asymp [\log(1/\delta_n)]^{-1/(\gamma_0 + \gamma_1)}$$

(iv) Finite admixtures: if G_0 has $k < \infty$ support points, k known, then we obtain a parametric rate:

$$\epsilon_{m,n} \asymp [\log(mn)/m]^{1/2} + [(\log n)^{1/2}/n^{1/4}]^{\alpha^*}$$

Optimal transport 00000000 Domain adaptation 00000000000000 Distance of latent structures

Quantifying abstract dependence

Geometric sparsity of support

- Sparse covering number: covering ε-balls that are separated by O(ε) in distance
- G₀ is supersparse with non-negative parameters (γ₀, γ₁), if its support admits sparse-covering number K(ε) ≤ [log(1/ε)]^{γ₀}, and the measure on such covering balls is at least g(ε) ≥ [log(1/ε)]^{-γ₁}.
 - let $\Theta = [0, 2]$, G_0 is supported on $S = \{1/2^k | k \in \mathbb{N}, k \ge 1\} \cup \{0\}$, and $G_0(\{1/2^k\}) \propto k^{-\gamma_1}$ for any $k \in \mathbb{N}$ and some $\gamma_1 > 1$; then G_0 is a supersparse measure with parameters $\gamma_0 = 1$ and γ_1 .

Domain adaptation 000000000000000 Distance of latent structures

Quantifying abstract dependence

Geometric sparsity of support

- Sparse covering number: covering ε-balls that are separated by O(ε) in distance
- G₀ is supersparse with non-negative parameters (γ₀, γ₁), if its support admits sparse-covering number K(ε) ≤ [log(1/ε)]^{γ₀}, and the measure on such covering balls is at least g(ε) ≥ [log(1/ε)]^{-γ₁}.
 - let $\Theta = [0, 2]$, G_0 is supported on $S = \{1/2^k | k \in \mathbb{N}, k \ge 1\} \cup \{0\}$, and $G_0(\{1/2^k\}) \propto k^{-\gamma_1}$ for any $k \in \mathbb{N}$ and some $\gamma_1 > 1$; then G_0 is a supersparse measure with parameters $\gamma_0 = 1$ and γ_1 .
- G_0 is ordinary sparse with parameters (γ_0, γ_1) if $K(\epsilon) \leq (1/\epsilon)^{\gamma_0}$, and $g(\epsilon) \geq \epsilon^{\gamma_1}$.

Distance of latent structures

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Geometric sparsity of support

- Sparse covering number: covering ε-balls that are separated by O(ε) in distance
- G₀ is supersparse with non-negative parameters (γ₀, γ₁), if its support admits sparse-covering number K(ε) ≤ [log(1/ε)]^{γ₀}, and the measure on such covering balls is at least g(ε) ≥ [log(1/ε)]^{-γ₁}.
 - let $\Theta = [0,2]$, G_0 is supported on $S = \{1/2^k | k \in \mathbb{N}, k \ge 1\} \cup \{0\}$, and $G_0(\{1/2^k\}) \propto k^{-\gamma_1}$ for any $k \in \mathbb{N}$ and some $\gamma_1 > 1$; then G_0 is a supersparse measure with parameters $\gamma_0 = 1$ and γ_1 .
- G_0 is ordinary sparse with parameters (γ_0, γ_1) if $K(\epsilon) \leq (1/\epsilon)^{\gamma_0}$, and $g(\epsilon) \geq \epsilon^{\gamma_1}$.
- Ordinary sparse measures are studied in fractal geometry: γ_0 is the Hausdorff dimension of the support, while γ_1 is the packing dimension.
 - if $\Theta = [0, 1]$, then the Hausdorff measure on the Cantor set is ordinary sparse with $\gamma_0 = \gamma_1 = \log 2/\log 3$.

Optimal transport 00000000 Domain adaptation 0000000000000 Distance of latent structures

Quantifying abstract dependence 00000000

Summarizing (this section)

- Optimal transport provides the first (and only so far) posterior contraction analysis of a hierarchical Bayesian model (Nguyen, Bernoulli 2016)
 - recent advances via harmonic analysis of mixture of product distributions (Wei & Nguyen, Annals 2022)
- It also connects to a notion of barycenter among a collection of measures (Agueh and Carlier, 2012), and gave rise to optimal transport based multi-level clustering methods (Ho et al, ICML 2017, Huynh et al, JMLR 2021)
- Full treatment of inverse bound for this general hierarchy remains elusive!
 - but there have been promising recent contributions from Nhat Ho, Aritra Guha, Yun Wei, Dat Do, Linh Do, Sunrit Chakraborty
 - for more details see my second lecture at VIASM on Thursday

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence •0000000

Outline

Data and distances

Optimal transport

Domain adaptation

Distance of latent structures Inverse bounds

Quantifying abstract dependence

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

Data and distances 0000	Optimal transport 00000000	Domain adaptation 00000000000000	Distance of latent structures 0000000 0000000000000000000000000000	Quantifying abstract dependence 0000000
----------------------------	-------------------------------	-------------------------------------	--	--

• <u>Story so far</u>: Optimal transport distance based inference on space of *distributions of data*, space of *distribution of quantities of interest*, distributions of *distributions*, and so on

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

• Other abstract notion of dependence can be quantified too:

- multivariate rank and quantiles
- independence
- exchangeability, partial exchangebility

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

Multivariate ranking

Given data sample $X_1, \ldots, X_n \in \mathbb{R}^d$, how to rank them?



Define empirical measure $\mu_n = \sum_{i=1}^n \delta_{X_i}$. Let ν_n be a discrete approximation of $\text{Unif}[0,1]^d$ (supported on a regular lattice). Then, the empirical rank function F is one which solves

$$F = \operatorname{argmin}_F \int \|x - F(x)\|^2 d\mu_n$$

such that $F_{\#}\mu_n = \nu_n$.

Deb and Sen (2019)
Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

Independence

Given observations of two random variables X and Y on Polish spaces \mathcal{X}, \mathcal{Y} , resp. Are X and Y independent or not?

Let
$$X \sim \mu; \ Y \sim
u;$$
 and $(X, Y) \sim \gamma$

- Shannon's mutual information: $I(X, Y) = K(\gamma, \mu \otimes \nu) = \int \log(d\gamma/d(\mu \otimes \nu))d\gamma$ provided $\gamma \ll \mu \otimes \nu$ (and $I(X, Y) = +\infty$ otherwise). Clearly I(X, Y) = 0 if and only if X and Y are independent.
- OT based dependence: (Nies et al, 2021; Wiesel, 2021)

$$\tau(X,Y) := T_c(\gamma,\mu\otimes\nu) = \inf_{\pi\in\Pi(\gamma,\mu\otimes\nu)} \int c((x,y),(x',y'))d\pi$$

$$\tau^Y(X,Y) := \int T_{c_Y}(\gamma_x,\nu)\mu(dx)$$

where c, c_Y is some notion of cost on $\mathcal{X} \times \mathcal{Y}$ and \mathcal{Y} , resp. γ_x denotes the conditional distribution of Y given X = x under the joint distribution γ (i.e., disintegration of γ wrt x)

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Optimal transport 00000000 Domain adaptation 00000000000000

Quantifying abstract dependence

Example: $\xi \sim \text{Unif}[0, 1]$ and $\zeta = f_n(\xi) \in \text{Unif}[0, 1]$ for zigzag functions f_n with *n* linear segments. n = 1 in (a), n = 8 in (b).

Note that $I(\xi, \zeta) = \infty$ in both cases for it fails to account for the metric structure of the support.



Illustration of Nies et al (2021)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

Distance to exchangeability

Given two exchangeable populations

$$X_1,\ldots,X_n|P_1 \stackrel{iid}{\sim} P_1; \quad Y_1,\ldots,Y_n|P_1 \stackrel{iid}{\sim} P_2$$

Can we tell whether $P_1 \perp P_2$, or $P_1 = P_2$ (almost surely)?

Catalano et al (2021):

- assume $(P_1, P_2) = (T(\mu_1), T(\mu_2))$, where $\mu_1 \stackrel{d}{=} \mu_2$
 - (μ_1,μ_2) random measures with joint independent increments
 - T is a generic map (e.g., normalization, exponential, kernel mixtures)
 - the above assumption covers a large range of model for two-sample problems
- the problem can be formulated as measuring the distance between complete random measures, which in turn boils down to the distance of the corresponding Lévy measures

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

Summary

- OT is a useful device for quantifying the space of <u>distributions of data</u>, <u>distribution of latent quantities</u> of interest, <u>distributions of distributions</u>, abstract notions of dependence
 - pushforward measures may be used as a useful way for modeling distributions, or to characterize convergence behavior of algorithms
- many open questions surrounding the mathematics behind the presence of latent-variable representation of complex models, and the <u>statistical</u> and <u>computational</u> theory for achieving and explicating the emergence of such representation
 - OT framework provides a promising approach

Optimal transport 00000000 Domain adaptation

Quantifying abstract dependence

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Selected References

- J. Zhu, A. Guha, D. Do, M. Xu, X. Nguyen and D. Zhao. Functional optimal transport: map estimation and domain adaptation for functional data. arXiv:2102.03895.
- X. Nguyen. Convergence of latent mixing measures in finite and infinite mixture models. Annals of Statistics, 2013.
- X. Nguyen. Borrowing strength in hierarchical Bayes: posterior concentration of the Dirichlet base measure. Bernoulli, 2016.
- Y. Wei and X. Nguyen. Convergence of de Finetti's mixing measures in latent structure models for exchangeable sequences. arXiv:2004.05542. Annals of Statistics, to appear.
- D. Do, N. Ho and X. Nguyen. Beyond black box densities: parameter learning for the deviated components. arXiv:2202.02651.