

Transport

Lecture 3



- GOALS //
- Define Transport Between Measures.
  - Properties of Transport
  - Variational Inference via Transport
  - Beyond KL Divergence.

WHAT IS TRANSPORT? //

Definition /  $\rho$  pdf on  $\mathbb{R}^d$ ;  $g: \mathbb{R}^d \rightarrow \mathbb{R}^c$ .

$\rho_g$  pdf, on  $\mathbb{R}^c$ , of  $g(z)$ ,  $z \sim \rho$ .

$\rho_g$  is the pushforward of  $\rho$  under  $g$ , written

$$\rho_g = g\# \rho.$$

Remark /  $A \subseteq \mathbb{R}^c$  Borel set then

$$\mathbb{P}^{\rho_g}(A) = \mathbb{P}^\rho(g^{-1}(A))$$

Definition / Let  $\rho_1, \rho_2 \in \mathcal{P}(\mathbb{R}^d)$ . A transport map  $g: \mathbb{R}^d \rightarrow \mathbb{R}^d$  between  $\rho_1$  &  $\rho_2$  is map with property  $\rho_2 = g\# \rho_1$ .

-

PROPERTIES OF TRANSPORT //

Important Remark / Recall the Wasserstein <sup>(L)</sup>  
metric in Kantorovich form

$$W_p(\rho_1, \rho_2) = \inf_{\pi \in \Pi_{\rho_1, \rho_2}} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|z - u\|^p \, \pi(z, u) dz du$$

Under certain smoothness assumptions this can be expressed in Monge form using transport:

$$W_p(\rho_1, \rho_2)^\phi = \inf_{\substack{g: g \# \rho_1 = \rho_2}} \int_{\mathbb{R}^d} \|z - g(z)\|_p^p \rho_1(z) dz$$

because  $\pi^*(z, u) = \delta(u - g^*(z)) \rho_1(z)$

$\pi$                            $g$

infimizing                          infimizing

Theorem / Let  $T \in C^1(\mathbb{R}^d, \mathbb{R}^d)$  be invertible. Then

$$\mathcal{D}_{\mathcal{K}_r}(\rho_1 \parallel \rho_2) = \mathcal{D}_{\mathcal{K}_r}(T_{\#}\rho_1 \parallel T_{\#}\rho_2).$$

Lemma /  $\pi \in \mathbb{R}^d$ ,  $q \in C(\mathbb{R}^d, \mathbb{R}^d)$  invertible. Then

$$\log q_{\#} \pi(u) = \log(\pi \circ q^{-1})(u) + \log \det D(q^{-1}(u))$$

(3)

Theorem / Let  $T \in C^1(\mathbb{R}^d, \mathbb{R}^d)$  be invertible. Then

$$D_{KL}(p_1 \parallel p_2) = D_{KL}(T_\# p_1 \parallel T_\# p_2).$$

Lemma /  $\pi \in \mathbb{R}^d$ ,  $q \in C^1(\mathbb{R}^d, \mathbb{R}^d)$  invertible. Then

$$\log q_\# \pi(u) = \log(\pi \circ q^{-1})(u) + \log \det D(q^{-1}(u))$$

## VARIATIONAL INFERENCE VIA TRANSPORT //

Recall Bayes Theorem:

$$\pi(u) = \frac{1}{Z} \delta(u) \rho(u), \quad Z = \mathbb{E}_{u \sim \rho} [\delta(u)].$$

We aim to find an approximate transport

$T(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with property  $T(\cdot; \theta)_\# \rho \approx \pi$ .

How should we choose  $\theta$ ? Define

$$F(\theta) = D_{KL}(T(\cdot; \theta)_\# \rho \parallel \pi).$$

By preceding theorem

$$F(\theta) = D_{KL}(\rho \parallel T^{-1}(\cdot; \theta)_\# \pi)$$

By properties of  $D_{KL}$

$$F(\theta) = -\mathbb{E}^{\text{unp}} \log(T^{-1}(\cdot; \theta) \# \pi) + \mathbb{E}^{\text{unp}} \log \rho(u)$$

(4)

Minimizing  $F(\cdot)$  is equivalent, by Lemma, to minimizing

$$-\mathbb{E}^{\text{unp}} [\log \pi_{\theta T}(u; \theta) + \log \det D_u T(u; \theta)]$$

$$\text{But } \log \pi_{\theta T} = \log \rho_T + \log \lambda_T - \log Z$$

& since  $Z$  is independent of  $T$ , and hence  $\theta$ , our minimization problem is equivalent to

$$\theta^* \in \arg \min_{\theta \in \Theta} J(\theta),$$

$$J(\theta) = -\mathbb{E}^{\text{unp}} [\log \rho_{\theta T}(u; \theta) + \log \lambda_{\theta T}(u; \theta) + \log \det D_u T(u; \theta)]$$

This is variational inference over set

$$Q = \{q = q = T \# \rho, T \in \mathcal{V}_{\Theta}\}$$

$\mathcal{V}$  is  $C^1$  invertible maps:

$$\mathcal{P} = \{ T(\cdot; \theta) \in C^1(\mathbb{R}^d, \mathbb{R}^d), \theta \in \mathbb{H} \mid$$

$$\det D_u T(u; \theta) > 0 \quad \forall (u, \theta) \in \mathbb{R}^d \times \mathbb{H} \}$$

(S)

## BEYOND KL DIVERGENCE //

Preceding requires explicit knowledge of  $\rho$ .  
 what if only samples from  $\rho$  are known,  
 not  $\rho$  itself:

Data Assumption // we are given

$$\{ u^{(n)} \}_{n=1, \dots, N} \text{ i.i.d. } \sim \rho$$

then define

$$\rho^N(u) = \frac{1}{N} \sum_{n=1}^N \delta_{u^{(n)}}(u)$$

$$\pi^N(u) = \frac{1}{Z} \varphi(u) \rho^N(u), \quad Z = \mathbb{E}^{u \sim \rho^N} [\varphi(u)].$$

It follows that

$$\pi^N(u) = \sum_{n=1}^N \omega^{(n)} \delta_{u^{(n)}}(u)$$

$$\omega^{(n)} = \varphi(u^{(n)}) / \sum_{m=1}^N \varphi(u^{(m)}), \quad \varphi(u^{(n)}) = \varphi(u^{(n)})$$

Find  $\theta$  to minimize

(6)

$$J^N(\theta) = D_E(T(\cdot; \theta), \#P^N, \pi^N).$$

This loss function can be evaluated  
using only samples from  $P^N$  &  $\pi^N$ .