

Optimal transport (OT) / Conditional transport(CT) and their applications

Dongsheng Wang

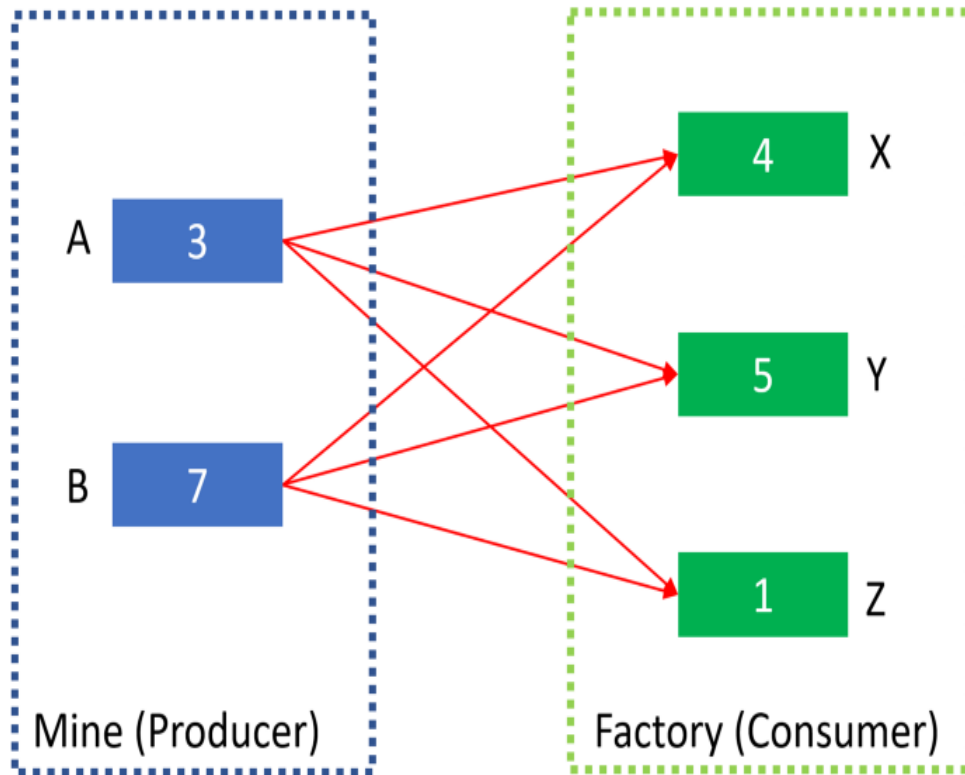
2022.10.12

➤ Outline

- ◆ Optimal transport theory
- ◆ Sinkhorn distances
- ◆ Conditional transport
- ◆ Applications and examples

➤ Optimal transport (OT)

◆ Let's start with an simple example



✓ The transport cost (pre-defined)

$$M \in R^{2 \times 3} \quad C_{ij} > 0$$

✓ The transport plan (need to-be-optimized)

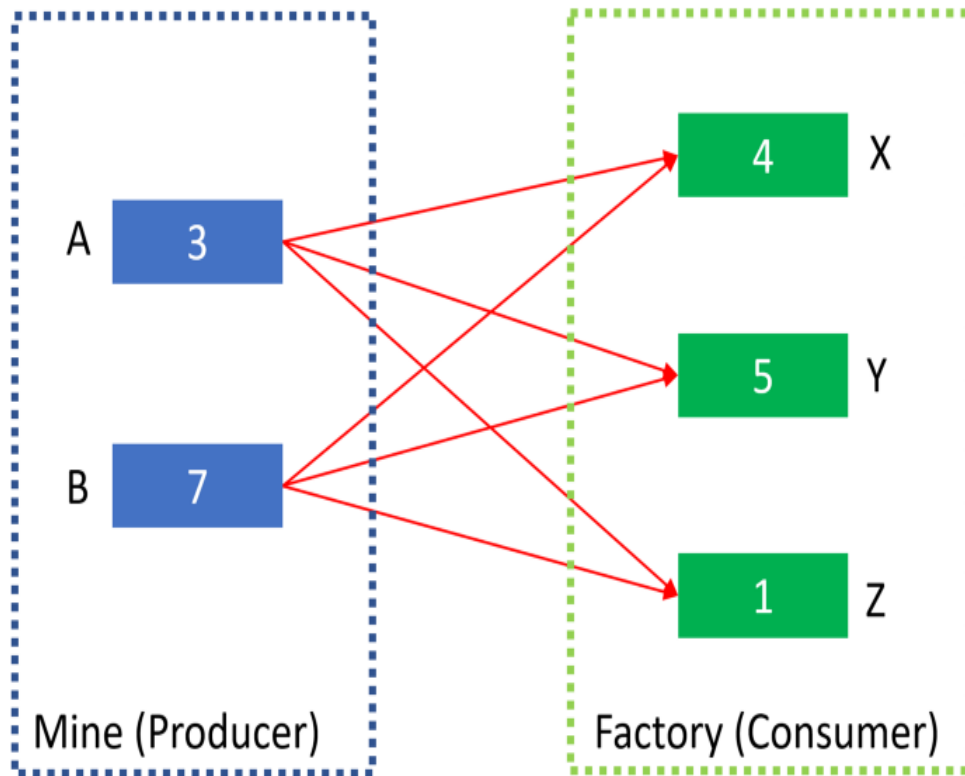
$$P \in R^{2 \times 3} \quad \sum_i P_{ij} = [4, 5, 1] \quad \sum_j P_{ij} = [3, 7]^T$$

✓ The total transport cost from the producer to consumer

$$L = \sum_{i \in \{A, B\}} \sum_{j \in \{X, Y, Z\}} P_{ij} M_{ij}$$

➤ Optimal transport (OT)

◆ Let's start with an simple example



$$d = \min \sum_{ij} P_{ij} C_{ij}$$

$$P_{ij} \geq 0$$

$$\sum_i P_{ij} = c_j$$

$$\sum_j P_{ij} = r_i$$

➤ Optimal transport (OT)

◆ The OT problem

Two empirical sets:

$$\{x_i\}_{i=1}^n \in \mathcal{X}^n \quad \{y_i\}_{i=1}^m \in \mathcal{Y}^m$$

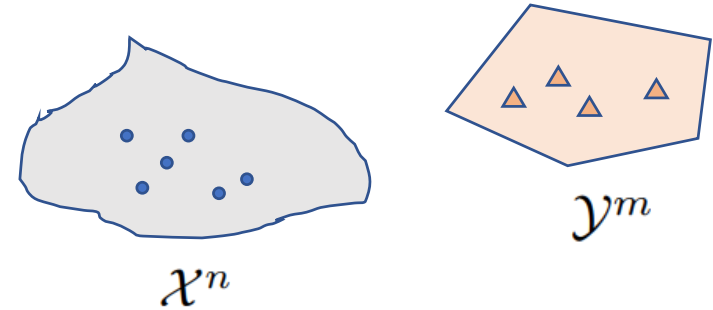
And the two corresponding weights over samples:

$$\mathbf{p} \in \mathbb{R}_+^n \text{ and } \mathbf{q} \in \mathbb{R}_+^m \text{ where } \sum_{i=1}^n \mathbf{p}_i = \sum_{i=1}^m \mathbf{q}_i = 1$$

The OT distance between \mathbf{p} and \mathbf{q} :

$$d_C(p, q) = \min_{\Gamma \in \Pi(\mathbf{p}, \mathbf{q})} \langle \Gamma, C \rangle, \quad \Pi(\mathbf{p}, \mathbf{q}) = \{\Gamma \in \mathbb{R}_+^{n \times m} \mid \gamma \mathbf{1}_m = \mathbf{p}, \gamma^T \mathbf{1}_n = \mathbf{q}\}$$

Where C is the transportation cost for each pair, e.g., $C_{ij} = \|x_i - y_j\|$



➤ Optimal transport (OT)

◆ The OT problem

Two empirical sets:

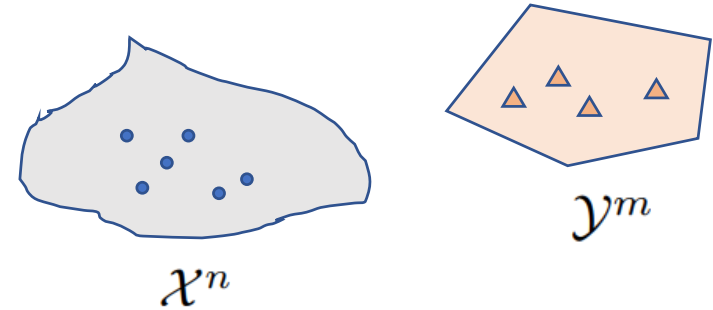
$$\{x_i\}_{i=1}^n \in \mathcal{X}^n \quad \{y_i\}_{i=1}^m \in \mathcal{Y}^m$$

And the two corresponding weights over samples:

$$\mathbf{p} \in \mathbb{R}_+^n \text{ and } \mathbf{q} \in \mathbb{R}_+^m \text{ where } \sum_{i=1}^n \mathbf{p}_i = \sum_{i=1}^m \mathbf{q}_i = 1$$

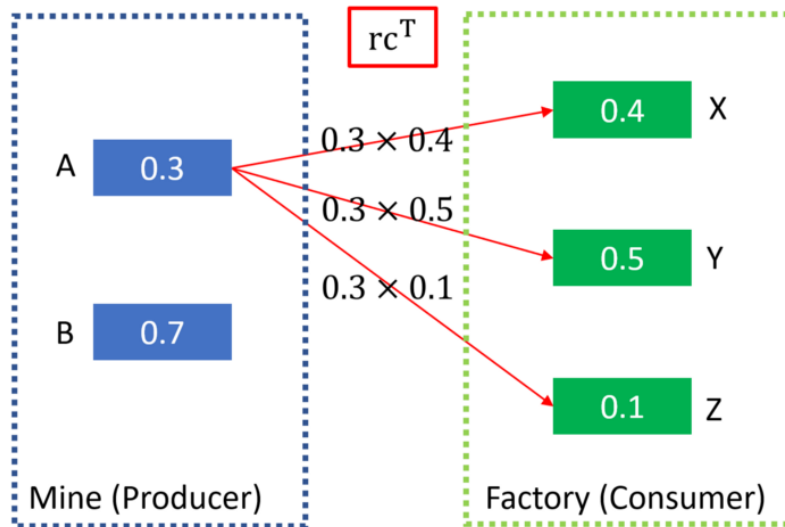
The OT distance between \mathbf{p} and \mathbf{q} :

$$d_C(p, q) = \min_{\Gamma \in \Pi(\mathbf{p}, \mathbf{q})} \langle \Gamma, C \rangle, \quad \Pi(\mathbf{p}, \mathbf{q}) = \{\Gamma \in \mathbb{R}_+^{n \times m} \mid \gamma \mathbf{1}_m = \mathbf{p}, \gamma^T \mathbf{1}_n = \mathbf{q}\}$$



The OT problem is a convex problem with the complexity $O(d^3 \log(d))$

➤ Sinkhorn distance



$$d = \min \sum_{i,j} P_{i,j} C_{i,j}$$

$$\begin{aligned} P_{i,j} &\geq 0 \\ \sum_j P_{i,j} &= r_i \\ \sum_i P_{i,j} &= c_j \end{aligned}$$



$$d^* = \min \sum_{i,j} P_{i,j} C_{i,j}$$

$$\begin{aligned} P_{i,j} &\geq 0 \\ \sum_j P_{i,j} &= r_i \\ \sum_i P_{i,j} &= c_j \end{aligned}$$

$$KL(P|rc^T) \leq \alpha$$

◆ What is the rc means

$$KL(P|rc^T) = \sum_{i,j} P_{i,j} \log \frac{P_{i,j}}{r_i c_j} = \sum_{i,j} P_{i,j} \log P_{i,j} - \sum_{i,j} P_{i,j} \log r_i - \sum_{i,j} P_{i,j} \log c_j$$

$$h(r) + h(c) - h(P) \leq \alpha$$

$$\hat{d} = \min \sum_{i,j} P_{i,j} C_{i,j} - \frac{1}{\lambda} h(P)$$

$$\begin{aligned} \sum_j P_{i,j} &= r_i \\ \sum_i P_{i,j} &= c_j \end{aligned}$$

➤ Sinkhorn distance

◆ Lagrange form of the Sinkhorn distance problem

$$L = \sum_{i,j} P_{i,j} C_{i,j} - \frac{1}{\lambda} h(P) + \sum_i m_i \left(\sum_j P_{i,j} - r_i \right) + \sum_j n_j \left(\sum_i P_{i,j} - c_j \right)$$

$$\frac{\partial L}{\partial P_{i,j}} = C_{i,j} + \frac{1}{\lambda} + \frac{1}{\lambda} \log P_{i,j} + m_i + n_j = 0$$

$$\begin{aligned} P_{i,j} &= e^{-\lambda m_i - 0.5} e^{-\lambda C_{i,j}} e^{-\lambda n_j - 0.5} \\ P_{i,j} &= u_i e^{-\lambda C_{i,j}} v_j \\ P &= \text{diag}(u) e^{-\lambda C} \text{diag}(v) \end{aligned}$$

➤ Sinkhorn algorithm

◆ The final objective

$$\min_{\Gamma \in \Pi(\mathbf{p}, \mathbf{q})} \langle \Gamma, C \rangle - \epsilon H(\Gamma), \text{ where } H(\Gamma) = - \sum_{i,j} \Gamma_{ij} \log \Gamma_{ij}$$

Algorithm 1 Computation of $d_M^\lambda(r, c)$ using Sinkhorn-Knopp's fixed point iteration

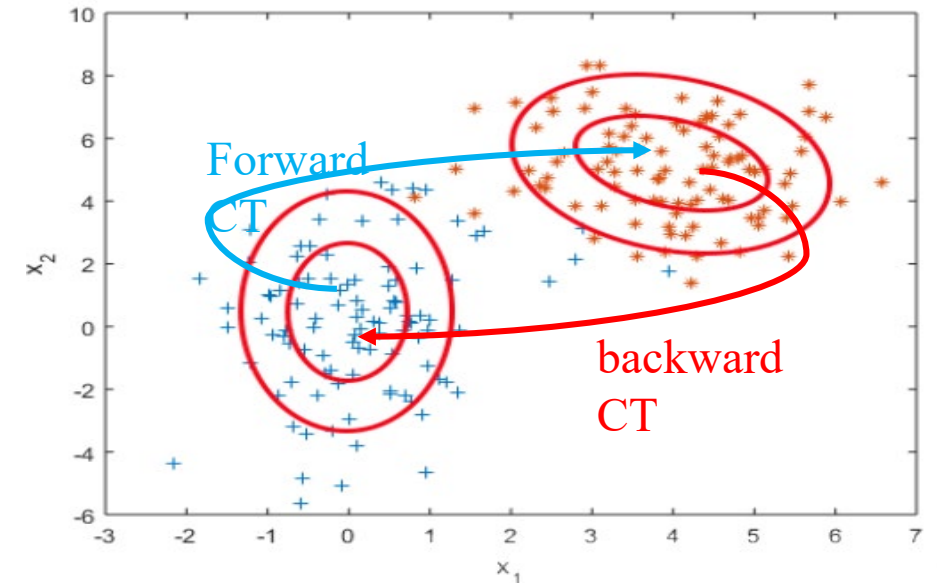
```
Input M,  $\lambda$ , r, c.  
I=(r>0); r=r(I); M=M(I,:); K=exp(- $\lambda$ *M)  
Set x=ones(length(r),size(c,2))/length(r);  
while x changes do  
    x=diag(1./r)*K*(c.*(1./(K'*(1./x))))  
end while  
u=1./x; v=c.*(1./(K'*u))  
 $d_M^\lambda(r,c)$ =sum(u.*(K.*M)*v)
```

➤ Conditional transport

◆ Drawbacks of OT

- The OT distance of p and q only considers single-direction transportation, e.g., from p to q
- There is an inner loop in sinkhorn algorithm which is inefficient.

◆ CT divergence is defined with a bidirectional distribution-to-distribution transport



➤ Conditional transport

◆ Forward CT is constructed in 3 steps

- Forward navigator

$$\pi(y | x) = \frac{e^{-d(x,y)} p_Y(y)}{\int e^{-d(x,y)} p_Y(y) dy}$$

Where, $d(x, y) = d(y, x)$ as a learnable distance function (e.g. $d(x, y) = \frac{(x - y)^2}{2e^\Phi}$)

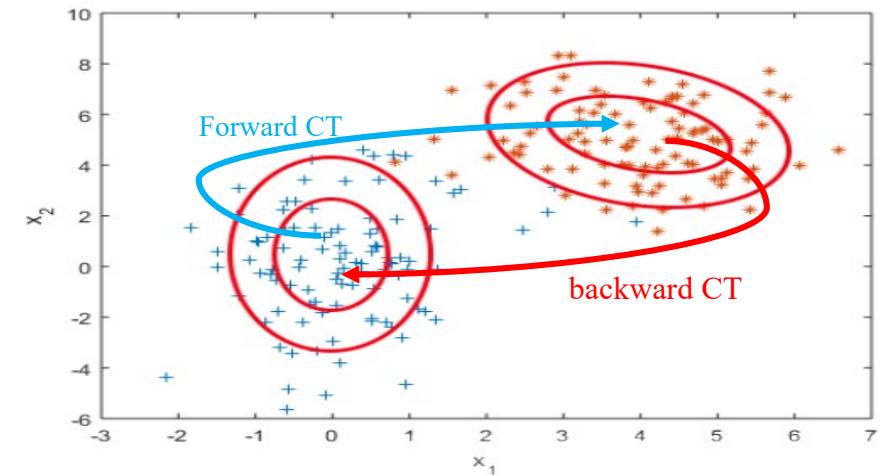
- Cost of forward transport for single point x

$$\text{cost}_x = \int c(x, y) \pi(y | x) dy$$

Where, $c(x, y) = c(y, x) \geq 0$ as the point-to-point transport cost (e.g. $c(x, y) = (x - y)^2$).

- Total cost of forward transport

$$\text{cost} = \int p_X(x) \int c(x, y) \pi(y | x) dx dy$$



➤ Conditional transport

✓ Backward CT in the same way

- backward navigator

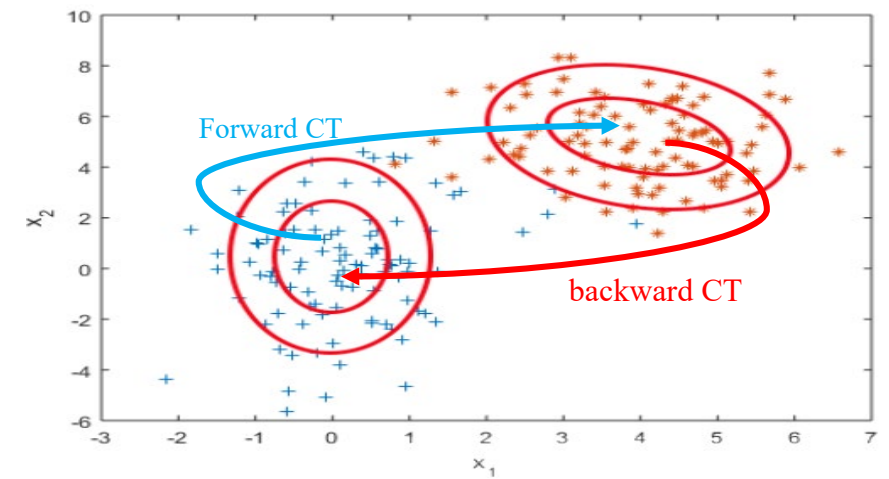
$$\pi(x | y) = \frac{e^{-d(x,y)} p_X(x)}{\int e^{-d(x,y)} p_X(x) dx}$$

- Cost of backward transport for single point x

$$\text{cost}_y = \int c(x, y) \pi(x | y) dx$$

- Total cost of backward transport

$$\text{cost} = \int p_Y(y) \int c(x, y) \pi(x | y) dx dy$$



➤ Conditional transport

- CT divergence

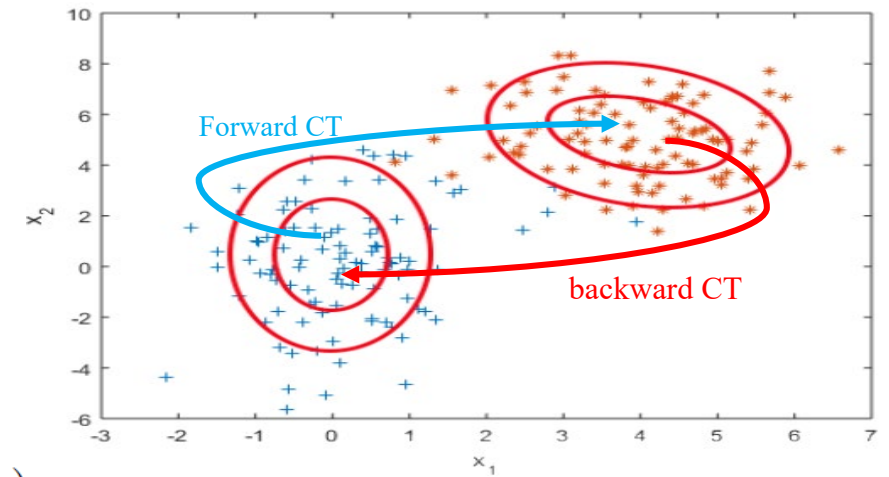
$$\mathcal{C}_{\phi, \theta}(\mu, \nu) \stackrel{\text{def.}}{=} \frac{1}{2} \mathcal{C}_{\phi, \theta}(\mu \rightarrow \nu) + \frac{1}{2} \mathcal{C}_{\phi, \theta}(\mu \leftarrow \nu)$$

Where, the forward CT and the backward CT:

$$\mathcal{C}_{\phi, \theta}(\mu \rightarrow \nu) = \mathbb{E}_{\mathbf{x} \sim p_X(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim \pi_{\phi}(\mathbf{y} | \mathbf{x})} [c(\mathbf{x}, \mathbf{y})], \quad \pi_{\phi}(\mathbf{y} | \mathbf{x}) \stackrel{\text{def.}}{=} \frac{e^{-d(\mathcal{T}_{\phi}(\mathbf{x}), \mathcal{T}_{\phi}(\mathbf{y}))} p_{\theta}(\mathbf{y})}{\int e^{-d(\mathcal{T}_{\phi}(\mathbf{x}), \mathcal{T}_{\phi}(\mathbf{y}))} p_{\theta}(\mathbf{y}) d\mathbf{y}},$$

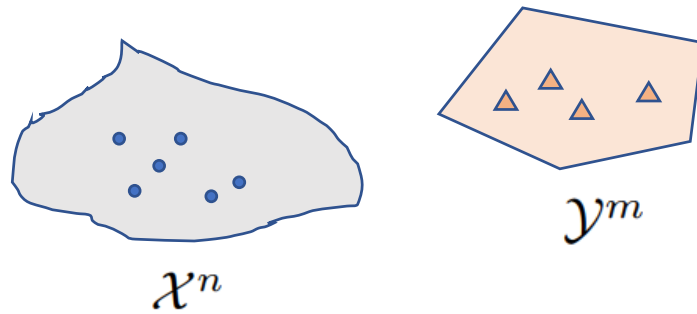
$$\mathcal{C}_{\phi, \theta}(\mu \leftarrow \nu) = \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y})} \mathbb{E}_{\mathbf{x} \sim \pi_{\phi}(\mathbf{x} | \mathbf{y})} [c(\mathbf{x}, \mathbf{y})], \quad \pi_{\phi}(\mathbf{x} | \mathbf{y}) \stackrel{\text{def.}}{=} \frac{e^{-d(\mathcal{T}_{\phi}(\mathbf{x}), \mathcal{T}_{\phi}(\mathbf{y}))} p_X(\mathbf{x})}{\int e^{-d(\mathcal{T}_{\phi}(\mathbf{x}), \mathcal{T}_{\phi}(\mathbf{y}))} p_X(\mathbf{x}) d\mathbf{x}},$$

$\mathcal{T}_{\phi}(\cdot) \in \mathbb{R}^H$, is a neural network based function.



➤ Applications

- ◆ Distance function: OT and CT that provide two convenient tools to measure distance between two empirical distributions (or two sets)
- ◆ The most core challenge is to design the P and Q according to your tasks



➤ Examples of topic model

- ◆ Two views of document from empirical distribution over word and topic space

$$p_j = \sum_v \tilde{x}_v \delta_{w_v} \quad q_j = \sum_k \tilde{\theta}_k \delta_{\beta_k}$$

$\tilde{x} \in \Delta^V$ Normalized Bag-of-Word vector $w_v \in R^d$ Embedding vector of v-th word

$\tilde{\theta} = f(\tilde{x}) \in \Delta^K$ Normalized topic proportion vector $\beta_k \in R^d$ Embedding vector of k-th topic

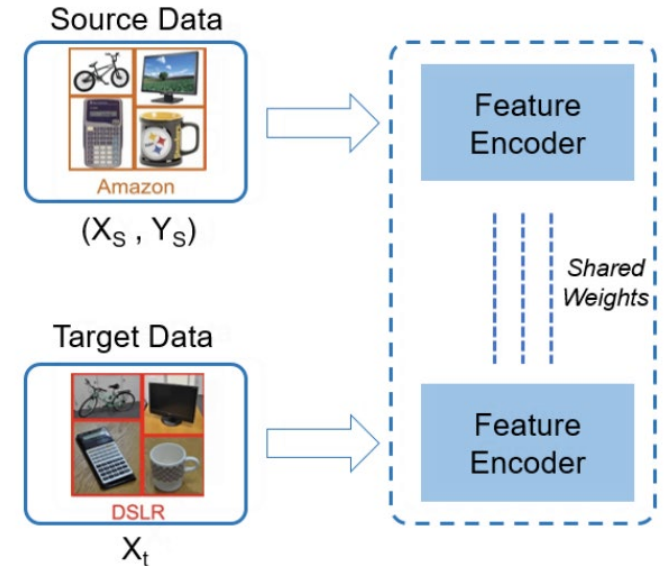
- ◆ p and q are two empirical distributions of the same document but with different supports (words and topics)

$$L = \sum_j d_C(p_j, q_j) \quad C_{vk} = 1 - \cos(w_v, \beta_k)$$

➤ Examples of domain adaptation

- ◆ The labeled source domain and the unlabeled target domain

$$\begin{aligned} \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{n_s} &\sim \mathcal{D}_s & \{\mathbf{x}_j^t\}_{j=1}^{n_t} &\sim \mathcal{D}_t^x \\ \mathbf{f}_i^s &= F_{\theta}(\mathbf{x}_i^s) & \mathbf{f}_j^t &= F_{\theta}(\mathbf{x}_j^t) \end{aligned}$$



- ◆ \mathcal{D}_s and \mathcal{D}_t have the same class label thus share similar class prototypes (centers) in the representation space

$$\mathcal{L}_{\text{cls}} = \mathbb{E}_{(\mathbf{x}_i^s, y_i^s) \sim \mathcal{D}_s} \left[\sum_{k=1}^K -\log p_{ik}^s \mathbf{1}_{\{y_i^s=k\}} \right], \quad p_{ik}^s := \frac{\exp(\boldsymbol{\mu}_k^T \mathbf{f}_i^s + b_k)}{\sum_{k'=1}^K \exp(\boldsymbol{\mu}_{k'}^T \mathbf{f}_i^s + b_{k'})}$$

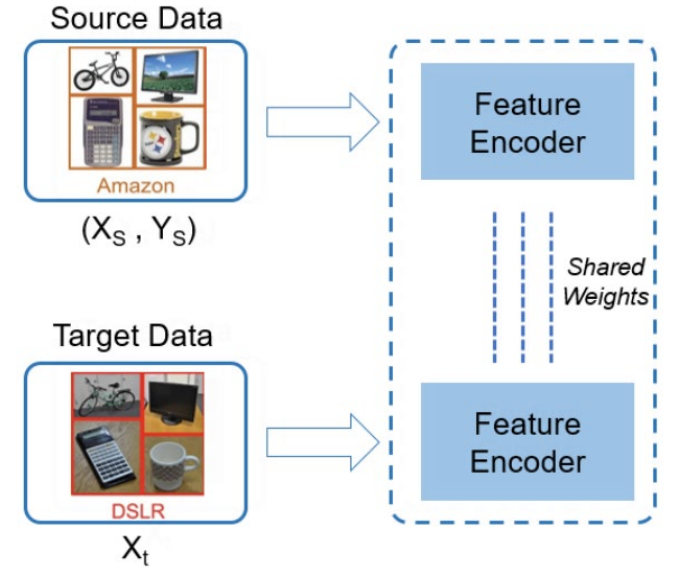
➤ Examples of domain adaptation

◆ Two set

$$P = \{f_j^t\}_{j=1}^M \quad Q = \{u_k\}_{k=1}^K$$

◆ CT is employed to finetune the encoder on the target set

$$CT = L_{t \rightarrow u} + L_{u \rightarrow t}$$



$$\mathcal{L}_{t \rightarrow \mu} = \mathbb{E}_{\mathbf{x}_j^t \sim \mathcal{D}_t^{\mathbf{x}}} \mathbb{E}_{\mu_k \sim \pi_{\theta}(\mu_k | \mathbf{f}_j^t)} [c(\mu_k, \mathbf{f}_j^t)] = \mathbb{E}_{\mathbf{x}_j^t \sim \mathcal{D}_t^{\mathbf{x}}} \left[\sum_{k=1}^K c(\mu_k, \mathbf{f}_j^t) \frac{p(\mu_k) \exp(\mu_k^T \mathbf{f}_j^t)}{\sum_{k'=1}^K p(\mu_{k'}) \exp(\mu_{k'}^T \mathbf{f}_j^t)} \right]$$

$$\begin{aligned} \mathcal{L}_{\mu \rightarrow t} &= \mathbb{E}_{\{\mathbf{x}_j^t\}_{j=1}^M \sim \mathcal{D}_t^{\mathbf{x}}} \mathbb{E}_{\mu_k \sim p(\mu_k)} \mathbb{E}_{\mathbf{f}_j^t \sim \pi_{\theta}(\mathbf{f}_j^t | \mu_k)} [c(\mu_k, \mathbf{f}_j^t)] \\ &= \mathbb{E}_{\{\mathbf{x}_j^t\}_{j=1}^M \sim \mathcal{D}_t^{\mathbf{x}}} \left[\sum_{k=1}^K p(\mu_k) \sum_{j=1}^M c(\mu_k, \mathbf{f}_j^t) \frac{\exp(\mu_k^T \mathbf{f}_j^t)}{\sum_{j'=1}^M \exp(\mu_k^T \mathbf{f}_{j'}^t)} \right] \end{aligned}$$

$$p(\mu_k)^{l+1} = \frac{1}{M} \sum_{j=1}^M \pi_{\theta}^l(\mu_k | \mathbf{f}_j^t), \quad \text{where} \quad \pi_{\theta}^l(\mu_k | \mathbf{f}_j^t) = \frac{p(\mu_k)^l \exp(\mu_k^T \mathbf{f}_j^t)}{\sum_{k'=1}^K p(\mu_{k'})^l \exp(\mu_{k'}^T \mathbf{f}_j^t)}$$

➤ Examples of imbalanced classification

◆ The imbalanced training data and the small balanced meta set

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \mathcal{D}_{\text{meta}} = \{(x_j, y_j)\}_{j=1}^M \quad M \ll N$$

◆ Conventional classification

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \ell(y_i, f(x_i; \boldsymbol{\theta}))$$

◆ The re-weight classification

$$\boldsymbol{\theta}^*(\boldsymbol{w}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N w_i l_i^{\text{train}}(\boldsymbol{\theta})$$

$$\boldsymbol{w}^* = \arg \min_{\boldsymbol{w}} \frac{1}{M} \sum_{j=1}^M l_j^{\text{meta}}(\boldsymbol{\theta}^*(\boldsymbol{w}))$$

➤ Examples of imbalanced classification

◆ Learn w from OT

$$P(\mathbf{w}) = \sum_{i=1}^N w_i \delta_{(x_i, y_i)^{\text{train}}} \quad Q = \sum_{j=1}^M \frac{1}{M} \delta_{(x_j, y_j)^{\text{meta}}}$$

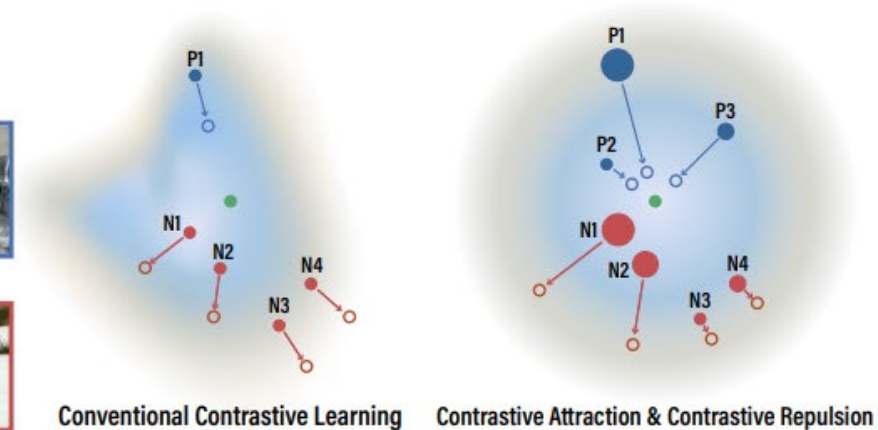
$$\min_{\mathbf{w}} \text{OT}(P(\mathbf{w}), Q) \stackrel{\text{def.}}{=} \min_{\mathbf{w}} \min_{\mathbf{T} \in \Pi(P(\mathbf{w}), Q)} \langle \mathbf{T}, \mathbf{C} \rangle$$

◆ The cost matrix

$$C_{ij} = d^{\text{Fea}}(\mathbf{z}_i^{\text{train}}, \mathbf{z}_j^{\text{meta}}) + d^{\text{Lab}}(y_i^{\text{train}}, y_j^{\text{meta}})$$

➤ Examples of contrastive learning

- ◆ The positive samples share close distance while the negative samples share far distance at the representation space



➤ Examples of contrastive learning

- ◆ Minimizing the total transport cost of mapping \mathbf{x} to its positive and negative sets.

$$L = E_{\mathbf{x} \sim p(\mathbf{x})} E_{\mathbf{x}^+ \sim \pi_{\theta}^+(\bullet | \mathbf{x})} [c(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^+))] + E_{\mathbf{x} \sim p(\mathbf{x})} E_{\mathbf{x}^- \sim \pi_{\theta}^-(\bullet | \mathbf{x})} [c(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^-))]$$

$$\pi_{\theta}^+(\mathbf{x}^+ | \mathbf{x}, \mathbf{x}_0) := \frac{e^{d_{t^+}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^+))} p(\mathbf{x}^+ | \mathbf{x}_0)}{Q^+(\mathbf{x} | \mathbf{x}_0)},$$

$$Q^+(\mathbf{x} | \mathbf{x}_0) =: \int e^{d_{t^+}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^+))} p(\mathbf{x}^+ | \mathbf{x}_0) d\mathbf{x}^+,$$

$$\pi_{\theta}^-(\mathbf{x}^- | \mathbf{x}) := \frac{e^{-d_{t^-}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^-))} p(\mathbf{x}^-)}{Q^-(\mathbf{x})},$$

$$Q^-(\mathbf{x}) := \int e^{-d_{t^-}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}^-))} p(\mathbf{x}^-) d\mathbf{x}^-,$$

$$c(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2;$$

$$d_{t^+}(\mathbf{x}, \mathbf{y}) = t^+ \|\mathbf{x} - \mathbf{y}\|_2^2, \quad t^+ \in \mathbb{R}_+;$$

$$d_{t^-}(\mathbf{x}, \mathbf{y}) = t^- \|\mathbf{x} - \mathbf{y}\|_2^2, \quad t^- \in \mathbb{R}_+.$$

➤ Reference

- [1] Zhao, He, et al. "Neural topic model via optimal transport." In ICLR2021
- [2] Huynh, Viet, He Zhao, and Dinh Phung. "OTLDA: A geometry-aware optimal transport approach for topic modeling." In NeurIPS2020.
- [3] Wang, Dongsheng, et al. "Representing Mixtures of Word Embeddings with Mixtures of Topic Embeddings." *ICLR2022*.
- [4] Tanwisuth, Korawat, et al. "A prototype-oriented framework for unsupervised domain adaptation." NeurIPS2021.
- [5] Guo, Dandan, et al. "Learning to Re-weight Examples with Optimal Transport for Imbalanced Classification." NeurIPS2022.
- [6] BERT-EMD: Many-to-Many Layer Mapping for BERT Compression with Earth Mover's Distance. In ACL2020.
- [7] **Learning Prototype-oriented Set Representations for Meta-Learning. In ICLR2022.**
- [8] Contrastive Attraction and Contrastive Repulsion for Representation Learning.
- [9] [Sinkhorn distances: Lightspeed computation of optimal transport](#). In NeurIPS2013.
- [9] From Word Embeddings To Document Distances. In ICML2015.

➤ Links

- [1] <https://amsword.medium.com/a-simple-introduction-on-sinkhorn-distances-d01a4ef4f085>
- [2] <https://towardsdatascience.com/optimal-transport-a-hidden-gem-that-empowers-todays-machine-learning-2609bbf67e59>
- [3] <http://alexhwilliams.info/itsneuronalblog/2020/10/09/optimal-transport/>
- [4] <https://medium.com/analytics-vidhya/introduction-to-optimal-transport-fd1816d51086>
- [5] <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6024256/>
- [6] <https://openreview.net/group?id=ICLR.cc/2023/Conference#all-submissions>