

COST PARAMETERIZATION AND CONSTRAINT RELAXATION FOR INVERSE OPTIMAL TRANSPORT WITH APPLICATIONS TO CONTRASTIVE LEARNING

Shanghai Jiao Tong University

April 25, 2023

IOT-CL

1	Contrastive Learning (CL)	2
1.1	Introduction	2
1.2	Loss Function	3
2	Optimal Transport (OT)	4
2.1	Introduction	4
2.2	Formalism	5
2.3	Regularization	6
2.4	Inverse Optimal Transport	7
2.5	IOT Inspires CL	8
3	Balanced Matching	10
3.1	Balanced Matching IOT-CL Loss	10
3.2	Gradient Epsilon	14
3.3	Alignment And Uniformity	15
3.4	Visualization	16
4	Appendix	18
4.1	Sinkhorn Algorithm	18

CONTRASTIVE LEARNING (CL)

INTRODUCTION

- ▶ Unsupervised Learning
- ▶ Transfer Knowledge
 - Pretrain + Downstream task
- ▶ Data Augmentation
 - Color Transformation
 - Geometric Transformation
 - Frame Order
- ▶ Architecture pipelines
 - End-to-End
 - Memory bank
 - Momentum Encoder
 - Clustering

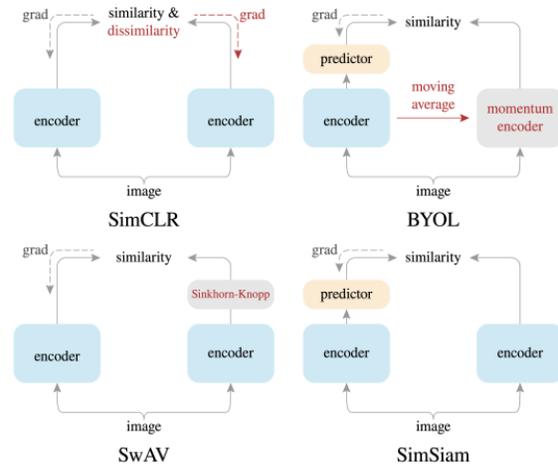


Figure. Comparison on CL architectures [textcitechen2021exploring].

CONTRASTIVE LEARNING (CL)

LOSS FUNCTION

Definition 1.1 (InfoNCE Loss)

The InfoNCE Loss Reads:

$$L_{\text{InfoNCE}} = \sum_{i=1}^n -\log \frac{\exp(s_{ii}/\tau)}{\exp(s_{ii}/\tau) + \sum_{k \neq i} \exp(s_{ik}/\tau)} \quad (1)$$

where $s_{ij} = \text{sim}(z_i, z'_j)$ is a similarity between the feature z_i and z'_j from 2 semantically related data.

Definition 1.2 (Alignment and Uniformity)

[Wang and Isola (2020)] view CL as enforcing 2 properties:

$$L_{\text{align}} = \sum_i \|z_i - z'_i\|_2^2 \quad \text{and} \quad L_{\text{uniform}} = \log \sum_{i,j} e^{2\|z_i - z'_j\|_2^2} \quad (2)$$

OPTIMAL TRANSPORT (OT)

INTRODUCTION

OPTIMAL TRANSPORT (OT)

INTRODUCTION

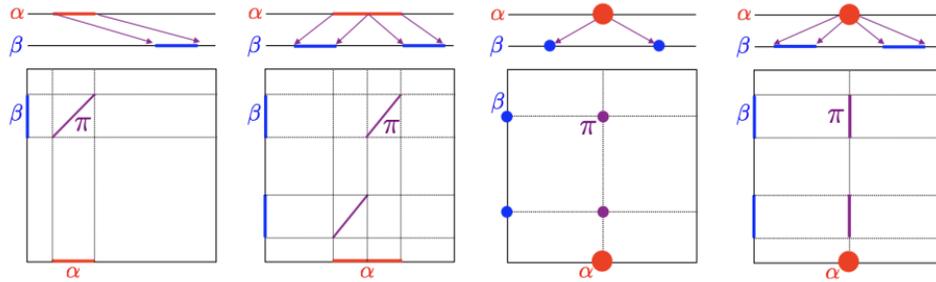


Figure. Four simple examples of optimal couplings between 1-D distributions, represented as maps above (arrows) and couplings below [From Peyré and Cuturi (2019)].

OPTIMAL TRANSPORT (OT)

FORMALISM

Definition 2.1 (Kantorovich's Optimal Transport Problem)

given the cost matrix \mathbf{C} , Kantorovich's OT involves solving the coupling \mathbf{P}

$$\min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \mathbf{P} \rangle = \sum_{i=1}^n \sum_{j=1}^m \mathbf{C}_{ij} \mathbf{P}_{ij} \quad (3)$$

where

$$U(\mathbf{a}, \mathbf{b}) = \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} \mid \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^T \mathbf{1}_n = \mathbf{b} \} \quad (4)$$

OPTIMAL TRANSPORT (OT)

FORMALISM

Definition 2.1 (Kantorovich's Optimal Transport Problem)

given the cost matrix \mathbf{C} , Kantorovich's OT involves solving the coupling \mathbf{P}

$$\min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \mathbf{P} \rangle = \sum_{i=1}^n \sum_{j=1}^m \mathbf{C}_{ij} \mathbf{P}_{ij} \quad (3)$$

where

$$U(\mathbf{a}, \mathbf{b}) = \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} \mid \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^T \mathbf{1}_n = \mathbf{b} \} \quad (4)$$

- ▶ When $n = m$ and $a = b = \mathbf{1}/n$, the OT is equivalent to solve a **balanced matching problem**.
- ▶ Similarly, we introduce some relaxation of $U(\mathbf{a}, \mathbf{b})$:

$$U(\mathbf{1}) = \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} \mid \mathbf{1}_n^T \mathbf{P} \mathbf{1}_m = 1 \} \quad \text{and} \quad U(\mathbf{a}) = \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} \mid \mathbf{P} \mathbf{1}_m = \mathbf{a} \}$$

OPTIMAL TRANSPORT (OT)

REGULARIZATION

Here we introduce **Regularization** which is highly suited to execution of GPU.

Definition 2.2

The objective reads:

$$\min_{\mathbf{P} \in U} \langle \mathbf{C}, \mathbf{P} \rangle - \epsilon H(\mathbf{P}) \quad (5)$$

where

$$H(\mathbf{P}) = - \sum_{i,j} \mathbf{P}_{ij} (\log \mathbf{P}_{ij} - 1). \quad (6)$$

- ▶ We call H Entropic Regularization. The function H is 1-strongly concave.
- ▶ Other form of Regularization: $G(\mathbf{P}) = \sum_{i,j} \left(-\frac{1}{2} \mathbf{P}_{ij}^2 + \mathbf{P}_{ij} \right)$. It's also 1-strongly concave.

OPTIMAL TRANSPORT (OT)

INVERSE OPTIMAL TRANSPORT

What if ...

- ▶ The cost matrix \mathbf{C} is unknown but the coupling \mathbf{P} is known?

OPTIMAL TRANSPORT (OT)

INVERSE OPTIMAL TRANSPORT

What if ...

- ▶ The cost matrix \mathbf{C} is unknown but the coupling \mathbf{P} is known?

Definition 2.3 (IOT problem)

By construct a min-min problem:

$$\min_{\theta} KL(\tilde{\mathbf{P}}|\mathbf{P}^{\theta}) \quad \text{where} \quad \mathbf{P}^{\theta} = \arg \min_{\mathbf{P} \in \mathcal{U}} \langle \mathbf{C}^{\theta}, \mathbf{P} \rangle - \epsilon H(\mathbf{P}) \quad (7)$$

where $\tilde{\mathbf{P}}$ is the ground truth and θ represents learnable parameters. Trivially, use Kullback–Leibler divergence to measure the distance between different distributions.

IOT INSPIRES CL

Example 2.1 (Equation 7 when $U = U(\mathbf{a})$)

The Lagrangian of the equation 7 reads:

$$L(\mathbf{P}, \lambda) = \langle \mathbf{C}^\theta, \mathbf{P} \rangle - \epsilon H(\mathbf{P}) - \sum_{i=1}^n \lambda_i \left(\sum_{j=1}^m \mathbf{P}_{ij} - \frac{1}{n} \right) \quad (8)$$

Through $\partial L / \partial P_{ij} = 0$, we have

$$\mathbf{P}_{ij}^\theta = \frac{\exp(-\mathbf{C}_{ij}^\theta / \epsilon)}{n \sum_{t=1}^m \exp(-\mathbf{C}_{it}^\theta / \epsilon)} \quad (9)$$

Setting $\tilde{\mathbf{P}} = \text{diag}(1/n, \dots, 1/n)$, the KL-divergence becomes:

$$KL(\tilde{\mathbf{P}} | \mathbf{P}^\theta) = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{\exp(-\mathbf{C}_{ii}^\theta / \epsilon)}{\sum_{t=1}^m \exp(-\mathbf{C}_{it}^\theta / \epsilon)} \right) + \text{Const} \quad (10)$$

IOT INSPIRES CL

Example 2.1 (Equation 7 when $U = U(\mathbf{a})$)

The Lagrangian of the equation 7 reads:

$$L(\mathbf{P}, \lambda) = \langle \mathbf{C}^\theta, \mathbf{P} \rangle - \epsilon H(\mathbf{P}) - \sum_{i=1}^n \lambda_i \left(\sum_{j=1}^m \mathbf{P}_{ij} - \frac{1}{n} \right) \quad (8)$$

Through $\partial L / \partial P_{ij} = 0$, we have

$$\mathbf{P}_{ij}^\theta = \frac{\exp(-\mathbf{C}_{ij}^\theta / \epsilon)}{n \sum_{t=1}^m \exp(-\mathbf{C}_{it}^\theta / \epsilon)} \quad (9)$$

Setting $\tilde{\mathbf{P}} = \text{diag}(1/n, \dots, 1/n)$, the KL-divergence becomes:

$$KL(\tilde{\mathbf{P}} | \mathbf{P}^\theta) = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{\exp(-\mathbf{C}_{ii}^\theta / \epsilon)}{\sum_{t=1}^m \exp(-\mathbf{C}_{it}^\theta / \epsilon)} \right) + \text{Const} \quad (10)$$

► which is [InfoNCE Loss!](#)

OPTIMAL TRANSPORT (OT)

IOT INSPIRES CL

Similarly, we could obtain the loss function under the circumstances in which $U = U(1)$.

$$L_H^{U(1)} = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{n \exp(-\mathbf{C}_{ii}^\theta / \epsilon)}{\sum_{t=-1}^n \sum_{s=1}^m \exp(-\mathbf{C}_{ts}^\theta / \epsilon)} \right) \quad (11)$$

However, if $U = U(\mathbf{a}, \mathbf{b})$, the closed-form coupling may not exist. We adopt [Sinkhorn algorithm](#) [Appendix 1] to approximate:

$$L_H^{U(\mathbf{a}, \mathbf{b})} = -\sum_{i=1}^m \sum_{j=1}^n \log \tilde{\mathbf{P}}_{ij} \left(\mathbf{P}_{ij}^\theta \right)^K \quad (12)$$

Use other regularization, like $G(\mathbf{P}) = \sum_{i,j} \left(-\frac{1}{2} \mathbf{P}_{ij}^2 + \mathbf{P}_{ij} \right)$

$$L_G^{U(\mathbf{a})} = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{n} \left(1 + \sum_{j=1}^n (\mathbf{C}_{ij}^\theta - \mathbf{C}_{ii}^\theta) \right) \right) \quad \text{and} \quad L_G^{U(1)} = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{n^2} \left(n + \sum_{s,t} (\mathbf{C}_{st}^\theta - \mathbf{C}_{ii}^\theta) \right) \right)$$

BALANCED MATCHING

BALANCED MATCHING IOT-CL LOSS

BALANCED MATCHING

BALANCED MATCHING IOT-CL LOSS

Given two collective points sets within $2N$ data points, we get features $\{z_i\}_{i=1}^N$ and $\{z_i^*\}_{i=1}^N$. We denote $\tilde{z}_{2k-1} = z_k, \tilde{z}_{2k} = z_k^*$.

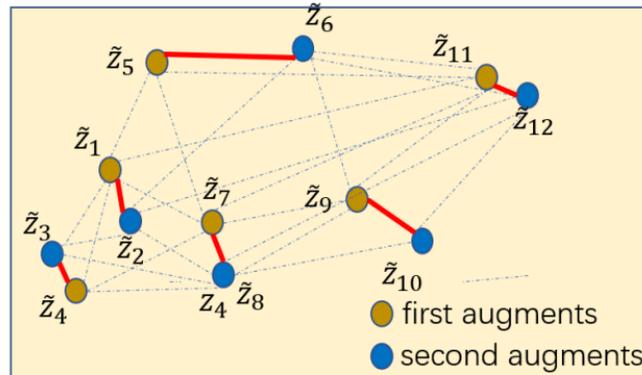


Figure. Balanced Matching

BALANCED MATCHING

BALANCED MATCHING IOT-CL LOSS

On the $2N$ data points, the cost matrix and $\tilde{\mathbf{P}}_{ij}$ are defined as follows:

$$\mathbf{C}_{ij}^{\theta} = \begin{cases} +\infty & \text{if } i = j \\ 1 - \tilde{s}_{ij} & \text{otherwise} \end{cases} \quad \text{and} \quad \tilde{\mathbf{P}}_{ij} = \begin{cases} 1/2N & \text{if } (i, j) \in S \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

where $S = S_1 \cup S_2$, S_1 and S_2 are defined the cost matrix and coupling as below:

$$S_1 = \{(i, j) \mid i = 2k, j = 2k - 1\} \quad S_2 = \{(i, j) \mid i = 2k - 1, j = 2k\} \quad \text{where } k = 1, \dots, N \quad (14)$$

BALANCED MATCHING

BALANCED MATCHING IOT-CL LOSS

Here we introduce different IOT-CL loss function for balanced matching problem under different constraint relaxations .

IOT-CL Loss Under U(a)

$$L_{\text{IOT-CL}}^{\text{U(a)}} = -\frac{1}{2N} \sum_{(i,j) \in \mathcal{S}} \log \left(\frac{\exp(-\mathbf{C}_{ij}^\theta / \epsilon)}{\sum_{s=1}^{2N} \mathbb{1}_{i \neq s} \exp(-\mathbf{C}_{is}^\theta / \epsilon)} \right) \quad (15)$$

This is the same as NT-Xent loss in SimCLR [Chen et al. (2020)].

IOT-CL Loss Under U(1)

$$L_{\text{IOT-CL}}^{\text{U(1)}} = -\frac{1}{2N} \sum_{(i,j) \in \mathcal{S}} \log \left(\frac{2N \exp(-\mathbf{C}_{ij}^\theta / \epsilon)}{\sum_{s=1}^{2N} \sum_{t=1}^{2N} \mathbb{1}_{s \neq t} \exp(-\mathbf{C}_{st}^\theta / \epsilon)} \right) \quad (16)$$

BALANCED MATCHING

BALANCED MATCHING IOT-CL LOSS

IOT-CL Loss Under $U(a,b)$

$$L_{\text{IOT-CL}}^{\text{U(a,b)}} = -\frac{1}{2N} \sum_{(i,j) \in S} \log(\mathbf{P}_{ij}^\theta)^K \quad (17)$$

where $(\mathbf{P}_{ij}^\theta)^K$ is solved by Sinkhorn algorithm, K is the iteration number.

IOT-CL Loss Under Gradient Constraint Relaxation

Under different constraint relaxations, \mathbf{P}_{ij}^θ will learn to approximate $\tilde{\mathbf{P}}_{ij}$ in different ways. We designed two gradient losses as below.

$$\begin{aligned} L_{\text{IOT-CL}}^{\text{Tighten}} &= L^{\text{U(a)}} \rightarrow L^{\text{U(a,b),K=1}} \rightarrow L^{\text{U(a,b),K=2}} \rightarrow L^{\text{U(a,b),K=4}} \rightarrow L^{\text{U(a,b),K=8}} \\ L_{\text{IOT-CL}}^{\text{Relax}} &= L^{\text{U(a,b),K=8}} \rightarrow L^{\text{U(a,b),K=4}} \rightarrow L^{\text{U(a,b),K=2}} \rightarrow L^{\text{U(a,b),K=1}} \rightarrow L^{\text{U(a)}} \end{aligned} \quad (18)$$

where the changing interval is fixed.

BALANCED MATCHING

GRADIENT EPSILON

IOT-CL Loss Under Gradient Epsilon

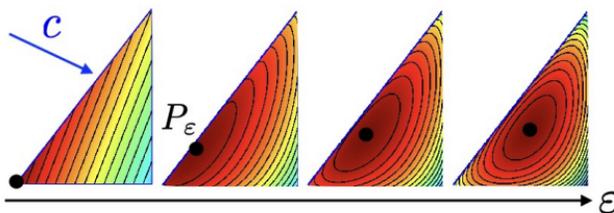


Figure. Entropic Regularization [Peyré and Cuturi (2019)].

Above figure illustrates the effect of the entropy to regularize a linear program over the simplex Σ_3 . We can see the entropy pushes the original LP solution away from the boundary of the triangle.

Thus the ϵ will control the "sharpness" of \mathbf{P}_{ij}^θ . A gradient setting will help \mathbf{P}_{ij}^θ move to $\tilde{\mathbf{P}}_{ij}$ faster.

BALANCED MATCHING

ALIGNMENT AND UNIFORMITY

We rethink the alignment and uniformity [Wang and Isola (2020)] from the matching prospective. By adding the uniformity penalty, we can get alignment and uniformity loss with matching view:

$$\min_{\theta} L_{\text{IOT-CL}}^{\text{Uniform}} = L_{\text{IOT-CL}} + \lambda_p \text{KL}(\bar{\mathbf{Q}}^{\theta} | \mathbf{P}^{\theta}) \quad (19)$$

where

$$\bar{\mathbf{Q}}_{ij}^{\theta} = \begin{cases} \mathbf{P}_{ij}^{\theta}, & \text{positive pair} \\ \text{mean}_{\text{negative pair}} \mathbf{P}_{ij}^{\theta}, & \text{negative pair} \end{cases} \quad (20)$$

The uniformity penalty will enhance our IOT-CL loss. The detailed experiment results will be available in our final paper.

BALANCED MATCHING

VISUALIZATION

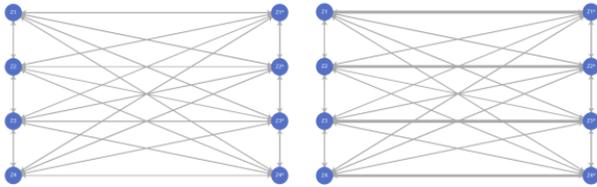


Figure. Visualization on balanced pair matching.

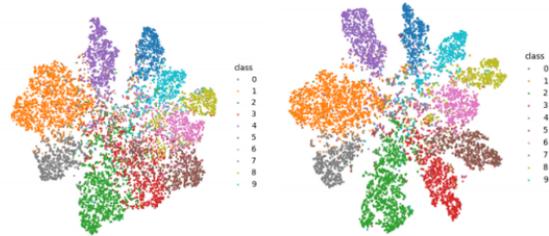


Figure. The effect of the Uniformity loss

REFERENCES I

-  [Chen, Ting et al. \(2020\)](#). “A Simple Framework for Contrastive Learning of Visual Representations”. In: *arXiv preprint arXiv:2002.05709*.
-  [Peyré, Gabriel and Marco Cuturi \(2019\)](#). “Computational Optimal Transport”. In: *Foundations and Trends in Machine Learning* 11.5-6, pp. 355–607.
-  [Wang, Tongzhou and Phillip Isola \(2020\)](#). “Understanding contrastive representation learning through alignment and uniformity on the hypersphere”. In: *International Conference on Machine Learning*. PMLR, pp. 9929–9939.

APPENDIX

SINKHORN ALGORITHM

Firstly, initialize \mathbf{P}^θ :

$$(\mathbf{P}^\theta)^0 = \exp(-\mathbf{C}^\theta/\epsilon) \quad (21)$$

Then update it step by step:

$$\begin{aligned} (\mathbf{P}^\theta)^k &\leftarrow \frac{1}{n} (\mathbf{P}^\theta)^{k-1} \oslash \left((\mathbf{P}^\theta)^{k-1} \mathbf{1}_{m \times m} \right) \\ (\mathbf{P}^\theta)^k &\leftarrow \frac{1}{m} (\mathbf{P}^\theta)^k \oslash \left(\mathbf{1}_{n \times n} (\mathbf{P}^\theta)^k \right) \end{aligned} \quad (22)$$

where the symbol \oslash represents element-wise division.