

The background of the slide features a repeating pattern of stylized, light pink flowers and leaves. The flowers have five petals and are connected by thin, curving stems with small leaves. The pattern is dense and covers the entire background.

Your contrastive learning problem is secretly a distribution alignment problem

Zihao Chen*, Chi-Heng Lin, Ran Liu, Jingyun Xiao, Eva L. Dyer
School of Electrical & Computer Engineering
Georgia Tech, Atlanta, GA

NeurIPS 2024

Outline

- ▶ Introduction: Contrastive Learning
 - ▶ Background
 - ▶ Motivation
- ▶ Generalized Contrastive Alignment (GCA)
 - ▶ Framework and Algorithm
 - ▶ Connections to Different CL Objectives
- ▶ Experiments
- ▶ Conclusion

Background

- ▶ Contrastive learning (CL) is a representation learning method that uses positive and negative pairs to define similarity in the latent space.
- ▶ Let $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ denote the dataset. For each sample \mathbf{x}_i in a batch of training data with size B , we create two augmented copies \mathbf{x}'_i and \mathbf{x}''_i independently. The $(\mathbf{x}'_i, \mathbf{x}''_i)$ is a positive pair, while $(\mathbf{x}'_i, \mathbf{x}''_j)$ for $j \neq i$ is a negative pair.
- ▶ **InfoNCE (INCE):** The loss function can be defined as :

$$\mathcal{L}_{\text{INCE}} = -\log \left(\frac{e^{s_{ii}}}{e^{s_{ii}} + \sum_{i \neq j} e^{s_{ij}}} \right) \quad (1)$$

where $s_{ij} = \varepsilon^{-1} f_{\theta}(\mathbf{x}'_i)^{\top} f_{\theta}(\mathbf{x}''_j) / \|f_{\theta}(\mathbf{x}'_i)\| \|f_{\theta}(\mathbf{x}''_j)\|$

Background

- ▶ **Robust InfoNCE (RINCE)** is a robust variant for InfoNCE with adjustable parameters λ and q to handle noisy:

$$\mathcal{L}_{\text{RINCE}}^{\lambda,q} = \frac{1}{q} \left(-e^{qs_{ii}} + \lambda^q \left(e^{s_{ii}} + \sum_{i \neq j} e^{s_{ij}} \right)^q \right) \quad (2)$$

- ▶ RINCE is a robust contrastive loss function characterized by its symmetric properties to noisy. ¹
 - ▶ When $q \rightarrow 1$, RINCE becomes a contrastive loss that fully satisfies the symmetry property (Ghosh et al. show that symmetric loss functions that is noise-tolerant ²)
 - ▶ When $q \rightarrow 0$, RINCE becomes asymptotically equivalent to InfoNCE.

¹Robust contrastive learning against noisy views (CVPR 2022)

²Making risk minimization tolerant to label noise (Neurocomputing 2015)

Background

► **Limitations of Traditional CL**

- CL's theoretical foundations and mechanisms for building representations are not fully understood.
- Traditional CL methods face challenges in handling real-world data, such as noise in views and the need for domain generalization.

► **Proposed Methods**

- Introduce the Generalized Contrastive Alignment (GCA) framework, which reinterprets CL as a distributional alignment problem.
- INCE and RINCE can be considered special cases under the GCA framework.
- Provide theoretical and experimental evidence of the benefits of GCA for noisy and generalized contrastive alignment.

GCA Framework

► Defining the Kernel Space:

- The augmentation kernel $K_\theta(\mathbf{x}'_i, \mathbf{x}''_j) = \exp\left(-\text{dist}(\tilde{f}_\theta(\mathbf{x}'_i), \tilde{f}_\theta(\mathbf{x}''_j))/\varepsilon\right)$ is defined, where $\text{dist}(\cdot)$ can be an arbitrary distance metric, and $\tilde{f}_\theta(\mathbf{x}'_i)$ is the normalized output of f_θ .

► Main Objective Formalization:

- The main objective can be defined as:

$$\min_{\theta} d_M(\mathbf{P}_{\text{tgt}} \parallel \mathbf{P}_{\theta}), \text{ with } \mathbf{P}_{\theta} = \arg \min_{\mathbf{P} \in \mathcal{B}} \{h(\mathbf{P}) + d_{\Gamma}(\mathbf{P} \parallel \mathbf{K}_{\theta})\} \quad (3)$$

- where $h(x)$ is a convex function (typically an indicator function),
- \mathcal{B} is a closed convex constraint set (e.g., Birkhoff polytope) 双随机矩阵多面体,
- d_{Γ} is a Bregman divergence,
- and d_M is a convex function (e.g., KL - divergence) that measures the divergence between \mathbf{P}_{θ} and \mathbf{P}_{tgt} .

GCA Framework-Advantages

- ▶ **Reframing CL as a Distribution Alignment Problem:**
 - ▶ Instead of just bringing positive pairs closer, GCA focuses on controlling the alignment of samples by defining a target transport plan \mathbf{P}_{tgt} .
 - ▶ For example, setting \mathbf{P}_{tgt} to resemble a diagonal matrix encourages each positive to align with itself or its augmentations, minimizing the deviation from an identity matrix (e.g., using KL - divergence).
- ▶ **Objective of GCA:**
 - ▶ Learn an encoder f_θ that minimizes the transport cost between positive samples. By defining \mathbf{P}_{tgt} with specific alignment rules, we can influence how samples are organized in the latent space.
 - ▶ This flexibility allows for encoding more nuanced forms of similarity and adapting to various learning tasks.

GCA Framework - Proximal Point Algorithm

Algorithm 1 Proximal-Point Algorithm for Generalized Contrastive Alignment (GCA)

- 1: **Initialization:** Initial encoder parameters θ , target transport plan \mathbf{P}_{tgt} , kernel function \mathbf{K}_θ , the function $h(x)$, divergences d_Γ and d_M (KL or W_1). Initialize transport plan \mathbf{P}_θ based on θ .
- 2: **Compute the transport coupling \mathbf{P}_θ :** Update \mathbf{P}_θ using the proximal operator scaling for fixed θ as described in Eq. (8):

$$\mathbf{P}_\theta = \arg \min_{\mathbf{P} \in \mathcal{B}} \{h(\mathbf{P}) + d_\Gamma(\mathbf{P} \parallel \mathbf{K}_\theta)\}.$$

- 3: **Calculate the loss:** Calculate deviation between the target and current transport plans

$$\mathcal{L}_{GCA} = d_M(\mathbf{P}_\theta, \mathbf{P}_{\text{tgt}}).$$

Update networks f_θ (encoder) and g_θ (projector) to minimize \mathcal{L}_{GCA} .

- 4: **Repeat until convergence:** Repeat steps 2 and 3 until convergence.
-

- Solve $d_\Gamma(\mathbf{P} \parallel \mathbf{K}_\theta)$ by Sinkhorn Algorithm based on iterative process;

Sinkhorn Algorithm

► Concept of Optimal Transport:

- Entropy-regularized OT (EOT) objective can be defined as :

$$\min_{\mathbf{P} \in \mathcal{B}} \langle \mathbf{P}, \mathbf{C} \rangle - \varepsilon H(\mathbf{P}) \quad (4)$$

- where $H(\mathbf{P}) = -\sum_{i,j} \mathbf{P}_{ij} \log(\mathbf{P}_{ij})$,
- Cost matrix: $C(x, y) = 1 - \langle x, y \rangle / \|x\| \|y\|$.
- Two constraints: $C_1^{(\mu)} := \{\mathbf{P} : \mathbf{P} \mathbf{1}_B = \mu\}$ and $C_2^{(\nu)} := \{\mathbf{P} : \mathbf{P}^\top \mathbf{1}_B = \nu\}$.

► Sinkhorn Algorithm and Bregman Projection:

- The first step finds the minimizer $\mathbf{P}^{(1)} = \arg \min \{ \epsilon KL(\mathbf{P} \| \mathbf{K}) : \mathbf{P} \mathbf{1}_B = \mu \}$ by the proximal operator $\text{Prox}_{C_1^{(\mu)}}^{KL}(\mathbf{K})$, compute its derivatives with respect to \mathbf{P} , set them to 0, and get \mathbf{P} , where $\mathbf{K}_{i,j} = \exp(-\mathbf{C}_{i,j}/\varepsilon)$.
- Subsequent steps project $\mathbf{P}^{(1)}$ onto the column constraint set.
 $\mathbf{P}^{(2)} := \text{Prox}_{C_2^{(\nu)}}^{KL}(\mathbf{P}^{(1)})$

Sinkhorn Algorithm

- The iterative updates can be succinctly expressed as the Sinkhorn iterations:

$$\mathbf{P}^{(2t+1)} = \text{diag}(\mathbf{u}^{(t+1)})\mathbf{K}\text{diag}(\mathbf{v}^{(t)}), \quad (5)$$

$$\mathbf{P}^{(2t+2)} = \text{diag}(\mathbf{u}^{(t+1)})\mathbf{K}\text{diag}(\mathbf{v}^{(t+1)}), \quad (6)$$

with the scaling vectors $\mathbf{u}^{(t)}$ and $\mathbf{v}^{(t)}$ updated according to:

$$\mathbf{u}^{(t+1)} \stackrel{\text{def}}{=} \frac{\mu}{\mathbf{K}\mathbf{v}^{(t)}}, \quad (7)$$

$$\mathbf{v}^{(t+1)} \stackrel{\text{def}}{=} \frac{\nu}{\mathbf{K}^T\mathbf{u}^{(t)}}. \quad (8)$$

GCA Framework - GCA - UOT Method

► Introduction to GCA - UOT:

- GCA - UOT is an extension of the GCA framework that utilizes unbalanced optimal transport (UOT) to relax the marginal constraints.

► Objective Function and Relaxation:

- The objective function is

$$\min_{\theta} d_M(\mathbf{P}_{\text{tgt}} \parallel \mathbf{P}_{\theta}) + \lambda_1 h_F(\mathbf{P}_{\theta} \mathbf{1} \parallel \mu) + \lambda_2 h_G(\mathbf{P}_{\theta}^{\top} \mathbf{1} \parallel \nu) + \varepsilon H(\mathbf{P}_{\theta}),$$

where h_F and h_G are divergence measures (e.g., KL divergence) that penalize deviations from the desired marginals μ and ν , and λ_1 and λ_2 are regularization parameters.

► Benefits and Impact:

- This relaxation allows the method to better handle data with imbalances or noise. It leads to different types of proximal operators and affects the coupling matrix. The impact of the entropy regularization parameter ε on the coupling matrix is studied, along with the number of iterations and corresponding sensitivity.

GCA Framework - Modifying the Target Transport Plan to Encode Matching Constraints

► Beyond Traditional Constraints

- In standard contrastive learning, the objective is often to minimize the deviations between the transport plan \mathbf{P}_θ and the identity matrix ($\mathbf{P}_{tgt} = I$). However, the GCA framework offers the flexibility to go beyond this simple one-to-one constraint.
- By modifying \mathbf{P}_{tgt} , we can incorporate additional meaningful constraints informed by domain knowledge or specific problem requirements.

GCA Framework - Modifying P_{tgt} to Encode Matching Constraints

► Domain Generalization Scenario

- Consider a domain generalization setting where each batch of data contains samples from multiple domains (e.g., Photo, Cartoon, Sketch, Art).
- We can structure the target alignment plan P_{tgt} as follows:

$$P_{tgt}[i,j] = I[i=j] + \alpha \cdot I(D_i = D_j, i \neq j) + \beta \cdot I(D_i \neq D_j, i \neq j)$$

Here, $I(\cdot)$ is the indicator function (equals 1 if the condition is true, 0 otherwise), D_i represents the domain of sample i , and $\alpha \geq 0$ and $\beta \geq 0$ are parameters.

- When $\alpha = 0$ and $\beta > 0$, we prioritize cross-domain alignment. This encourages the model to learn similarities and relationships between samples from different domains.
- Conversely, when $\alpha > 0$ and $\beta = 0$, the focus is on intra-domain alignment, emphasizing the similarities within each domain.

Connections to Different CL Objectives - Connection to INCE

► INCE as a Single Step in GCA

- An interesting connection exists between the GCA framework and the widely used InfoNCE (INCE) objective. We can interpret INCE as a single step within an iterative GCA objective.
- Under specific conditions (such as using the augmentation kernel \mathbf{K}_θ with cosine similarity, setting d_Γ and d_M to KL - divergence, and using the constraint set $\mathcal{C}_1^{(\mu)}$), the INCE objective can be re - expressed as a GCA problem.
- Specifically, the INCE objective in Equation (1) can be written as
$$\min_{\theta} KL \left(\mathbf{I} \parallel \text{Prox}_{\mathcal{C}_1^{(\mu)}}^{\mathbf{K}_\theta} (\mathbf{K}_\theta) \right).$$
 This shows that INCE can be seen as solving the matching problem with row normalization constraints $\mathcal{C}_1^{(\mu)}$.

Connections to Different CL Objectives

Table 1: *Comparison of different contrastive alignment objectives.* Here we have C_1^μ and C_2^ν as constraint sets (denoted as \mathcal{B}) defined in Equation (4) with their corresponding indicator function. "Iter" refers to iterative methods.

Methods	d_M	d_Γ	\mathcal{B}	Iter
INCE	KL	KL	C_1^μ	
GCA-INCE	KL	KL	$C_1^\mu \cap C_2^\nu$	✓
RINCE (q=1)	W1	KL	C_1^μ	
GCA-RINCE (q=1)	W1	KL	$C_1^\mu \cap C_2^\nu$	✓
BYOL	KL	L2	$R^{B \times B}$	

- Modification of the different parts of the main objective $(d_\Gamma, d_M, \mathcal{B})$ can be connected to different contrastive losses.
- GCA-INCE and GCA-RINCE are iterative processes.
- INCE is a single step in GCA-INCE. RINCE is a single step in GCA-RINCE.

Experiments

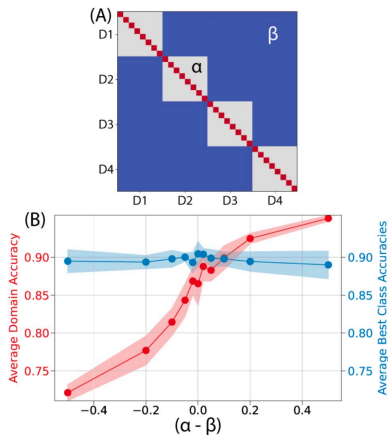
- ▶ Methods:
 - ▶ InfoNCE(**INCE**), Robust InfoNCE (**IRINCE**),
 - ▶ the GCA - based alternatives (**IGCA-INCE**, **IGCA-RINCE**)
 - ▶ **ISimCLR**, **IBYOL**, and **IOT**
 - ▶ **GCA-UOT**
- ▶ Datasets:
 - ▶ SVHN (Resnet-50)
 - ▶ ImageNet-100, CIFAR-10, CIFAR-100(Resnet-18)
 - ▶ The extreme angmentations datastes: CIFAR-10(Ex), CIFAR-100(Ex), CIFAR-10C(contains various corruptions), CIFAR-10C(Ex) (Resnet-18)

Experiments

Method	Standard Setting				Noisy Setting			
	CIFAR-10	CIFAR-100	SVHN	ImageNet100	CIFAR-10 (Ex)	CIFAR-100 (Ex)	CIFAR-10C	CIFAR-10C (Ex)
INCE	92.01 \pm 0.40	70.07 \pm 0.42	90.60 \pm 0.17	73.01 \pm 0.61	82.03 \pm 0.32	54.70 \pm 0.43	87.20 \pm 0.37	74.84 \pm 0.21
GCA-INCE	<u>92.36 \pm 0.24</u>	<u>70.11 \pm 0.45</u>	90.40 \pm 0.16	73.04 \pm 0.76	82.18 \pm 0.69	54.91 \pm 0.56	87.34 \pm 0.34	76.00 \pm 0.17
Δ	+0.35	+0.04	-0.20	+0.03	+0.15	+0.21	+0.14	+1.16
RINCE	91.05 \pm 0.50	69.06 \pm 0.64	90.97 \pm 0.19	71.91 \pm 0.43	82.60 \pm 0.63	55.43 \pm 0.48	88.62 \pm 1.33	77.05 \pm 0.82
GCA-RINCE	92.09 \pm 0.22	69.72 \pm 0.27	<u>91.45 \pm 0.41</u>	<u>73.44 \pm 0.55</u>	<u>82.76 \pm 0.49</u>	<u>55.90 \pm 0.41</u>	<u>88.76 \pm 0.72</u>	<u>77.23 \pm 0.76</u>
Δ	+1.04	+0.66	+0.48	+1.53	+0.16	+0.47	+0.14	+0.18
SimCLR	92.16 \pm 0.16	69.95 \pm 0.14	90.24 \pm 0.24	72.20 \pm 0.78	81.87 \pm 0.53	54.54 \pm 0.79	86.98 \pm 1.59	73.79 \pm 0.32
BYOL	90.56 \pm 0.59	69.75 \pm 0.37	89.50 \pm 0.46	69.75 \pm 0.83	81.55 \pm 0.50	54.18 \pm 0.46	87.88 \pm 1.02	69.40 \pm 1.11
IOT [51]	90.99 \pm 0.54	67.19 \pm 0.21	90.15 \pm 0.21	72.27 \pm 0.53	80.59 \pm 0.64	52.40 \pm 0.48	67.36 \pm 1.97	58.75 \pm 1.96
IOT-uni [51]	90.89 \pm 0.57	67.03 \pm 0.40	90.54 \pm 0.20	72.88 \pm 0.71	80.79 \pm 0.24	53.04 \pm 0.52	69.58 \pm 1.25	59.05 \pm 1.86
GCA-UOT	92.61 \pm 0.32	71.45 \pm 0.37	91.96 \pm 0.15	74.09 \pm 0.40	83.18 \pm 0.44	56.30 \pm 0.51	89.61 \pm 0.30	77.60 \pm 0.54

- ▶ The GCA- versions of INCE and RINCE exhibited performance gains in most settings.
- ▶ GCA - UOT achieved the top performance across all four datasets tested.
- ▶ The GCA-based strategies effectively enhanced the model's generalization ability and adaptability to aggressive data augmentations, especially in handling complex and noisy data.

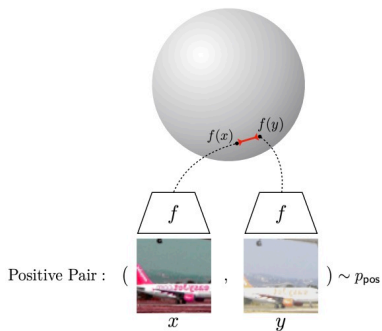
Experiments



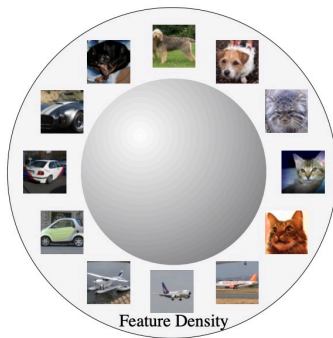
- ▶ Domain generalization task, samples originate from different domains (e.g. Photo, Cartoon, Sketch, Arts).
- ▶ The training was conducted on the PACS dataset using a ResNet18 encoder with the GCA-INCE objective.
- ▶ For \mathbf{P}_{tgt} , we set $\{\alpha = 0, \beta > 0\}$ to prioritize cross-domain alignment and $\{\alpha > 0, \beta = 0\}$ to focus on intra-domain alignment;
- ▶ Revealed that increasing the domain alignment weight enhances the accuracy of domain classification (from 72.11% to 95.16%) without diminishing classification performance.

Conclusion

- ▶ In this work, we introduced the Generalized Contrastive Alignment (GCA) framework, which reinterprets contrastive learning as a distribution alignment problem using optimal transport.
- ▶ We demonstrated the connections between GCA and different contrastive learning objectives, such as INCE, RINCE, and BYOL. Through theoretical analysis and experimental validation, we showed that GCA methods can **improve alignment and uniformity in the latent space**, leading to more discriminative and resilient representations.
- ▶ **Theoretical Analysis:**
 - ▶ Improved alignment with GCA;
 - ▶ Improved Uniformity of Representations Through GCA;
 - ▶ Impacts of GCA on a downstream classification task;



Alignment: Similar samples have similar features.
 (Figure inspired by [Tian et al. \(2019\)](#).)



Uniformity: Preserve maximal information.