



# How to Exploit Hyperspherical Embeddings for Out-of-Distribution Detection?

Yifei Ming<sup>1</sup>, Yiyou Sun<sup>1</sup>, Ousmane Dia<sup>2</sup>, Yixuan Li<sup>1</sup> Department of Computer Sciences, University of Wisconsin-Madison<sup>1</sup> Meta<sup>2</sup> {alvinming, sunyiyou, sharonli}@cs.wisc.edu, ousamdia@meta.com

ICLR 2023

# Background



#### The task of OOD Detection



Mathematically, let  $\mathcal{D}_{test}^{ood}$  denote an OOD test set where the label space  $\mathcal{Y}^{ood} \cap \mathcal{Y}^{in} = \emptyset$ . The decision can be made via a level set estimation:

$$G_{\lambda}(x) = \mathbb{1}\{S(x) \ge \lambda\},\$$

where samples with higher scores S(x) are classified as ID and vice versa. The threshold  $\lambda$  is typically chosen so that a high fraction of ID data (*e.g.* 95%) is correctly classified.

➤ Using the model's confidence score for OOD detection, which can be abnormally high on OOD samples.

Distance-based methods leverage feature embeddings extracted from a model, and operate under the assumption that the test OOD samples are relatively far away from the clusters of ID data.

# Background



# Hyperspherical embeddings



# Naturally modeled by the von Mises-Fisher (vMF) distribution:

$$p_d\left(\mathbf{z};\boldsymbol{\mu}_c,\boldsymbol{\kappa}\right) = \frac{\kappa^{\frac{p}{2}-1}}{(2\pi)^{\frac{p}{2}}I_{\frac{p}{2}-1}(\boldsymbol{\kappa})} e^{\kappa\vec{\mu}^T\vec{z}} = Z_d(\boldsymbol{\kappa})\exp\left(\kappa\boldsymbol{\mu}_c^{\top}\mathbf{z}\right)$$

A hypersphere is a topological space that is homeomorphic to a standard n-sphere, which is the set of points in (n + 1)-dimensional Euclidean space that are located at a constant distance from the center. When the sphere has a unit radius, it is called the unit hypersphere. Formally, an n-dimensional unit-hypersphere

$$S^{n} := \{ z \in \mathbb{R}^{n+1} | ||z||_{2} = 1 \}$$

Geometrically, hyperspherical embeddings lie on the surface of a hypersphere.

# **Motivation**



#### **Previous work**

> Arguably, the efficacy of distance-based approaches can depend largely on the quality of feature embeddings.

Recent works including SSD+ and KNN+ directly employ off-the-shelf contrastive losses for OOD detection.

They use the SupCon for learning the embeddings, which are then used for OOD detection with either parametric Mahalanobis distance or non-parametric KNN distance.

(Example) When trained on CIFAR-10 using SupCon loss, the average angular distance between ID and OOD data is only 29.86 degrees in the embedding space, which is too small for effective ID-OOD separation.

> Two properties

- Each sample has a higher probability assigned to the correct class in comparison to incorrect classes
- > Different classes are far apart from each other.

Method



#### **Framework Overview**

#### Compactness and Dispersion Regularized (CIDER)



A deep neural network encoder  $f : \mathcal{X} \mapsto \mathbb{R}^e$  that maps the augmented input to a high dimensional feature

A projection head:  $h : \mathbb{R}^e \to \mathbb{R}^d$  maps high-dim feature to a lower-dim feature





#### **Desirable Hyperspherical Embeddings in the Open World**



ID and OOD Samples on the hypersphere





#### Preliminaries

≻Multi-class classification

 $\mathcal{X}$  denotes the input space

 $\mathcal{Y}^{in} = \{1, 2, ..., C\}$  denotes the ID labels.

The training set  $\mathcal{D}_{tr}^{in} = \{(x_i, y_i)\}_{i=1}^N$  is drawn *i.i.d.* from  $P_{\mathcal{X}\mathcal{Y}^{in}}$ .

 $P_{\mathcal{X}}$  denote the marginal distribution over  $\mathcal{X}$ , which is called the in-distribution (ID).





#### **Modeling Hyperspherical Embeddings**

#### ► Intra-class compactness

- > Augmented input  $\tilde{x} \to \text{high-dim}$  feature  $f(\tilde{x}) \to \text{low-dim}$  feature  $\tilde{z} = h(f(\tilde{x}))$
- → Hyperspherical (  $l_2$ -normalized) feature  $z = \tilde{z}/\|\tilde{z}\|_2$

#### > von-Mises-Fisher (vMF) distribution

$$p_{d}(\mathbf{z};\boldsymbol{\mu}_{c},\kappa) = \frac{\kappa^{\frac{p}{2}-1}}{(2\pi)^{\frac{p}{2}}I_{\frac{p}{2}-1}(\kappa)} e^{\kappa \vec{\mu}^{T}\vec{z}} = Z_{d}(\kappa) \exp\left(\kappa \boldsymbol{\mu}_{c}^{\top}\mathbf{z}\right) \quad \mathbb{P}\left(y = c \mid \mathbf{z}; \left\{\kappa, \boldsymbol{\mu}_{j}\right\}_{j=1}^{C}\right) = \frac{Z_{d}(\kappa) \exp\left(\kappa \boldsymbol{\mu}_{c}^{\top}\mathbf{z}\right)}{\sum_{j=1}^{C}Z_{d}(\kappa) \exp\left(\kappa \boldsymbol{\mu}_{j}^{\top}\mathbf{z}\right)}$$
  

$$\succ \mu_{c} \text{ is the class prototype (mean) of class c} = \frac{\exp\left(\boldsymbol{\mu}_{c}^{\top}\mathbf{z}/\tau\right)}{\sum_{j=1}^{C}\exp\left(\boldsymbol{\mu}_{c}^{\top}\mathbf{z}/\tau\right)},$$

 $\succ \kappa$  indicates the spread of the distribution around mean direction





#### **Modeling Hyperspherical Embeddings**

#### ➢Inter-class dispersion

Optimize large angular distances among different class prototypes

 $\mathcal{L}_{\text{dis}} = \frac{1}{C} \sum_{i=1}^{C} \log \frac{1}{C-1} \sum_{j=1}^{C} \mathbb{1}\{j \neq i\} e^{\boldsymbol{\mu}_{i}^{\top} \boldsymbol{\mu}_{j}/\tau}.$ 

Compactness and Dispersion Regularized Learning (CIDER)

$$\mathcal{L}_{ ext{CIDER}} = \mathcal{L}_{ ext{dis}} + \lambda_c \mathcal{L}_{ ext{comp}}$$

➢Update of class prototypes

> Class prototypes are initialized as the classwise mean  $\mu_c$  using the training set

➤ Updated via exponential-moving-average (EMA):

$$\boldsymbol{\mu}_{c} := \text{Normalize} \left( \alpha \boldsymbol{\mu}_{c} + (1 - \alpha) \mathbf{z} \right), \forall c \in \{1, 2, \dots, C\}$$





#### **Pseudo-code for CIDER**

Algorithm 1: Pseudo-code of CIDER.

1 Input: Training dataset  $\mathcal{D}$ , neural network encoder f, projection head h, classifier q, class prototypes  $\mu_i$  $(1 \le j \le C)$ , weights of loss terms  $\lambda_d$ , and  $\lambda_c$ , temperature  $\tau$ **2** for epoch = 1, 2, ..., dofor iter = 1, 2, ..., do3 sample a mini-batch  $B = {\mathbf{x}_i, y_i}_{i=1}^{b}$ 4 obtain augmented batch  $\tilde{B} = {\{\tilde{\mathbf{x}}_i, \tilde{y}_i\}_{i=1}^{2b}}$  by applying two random augmentations to  $\mathbf{x}_i \in B$ 5  $\forall i \in \{1, 2, \dots, b\}$ for  $\tilde{\mathbf{x}}_i \in B$  do 6 // obtain normalized embedding  $\tilde{\mathbf{z}}_i = h(f(\tilde{\mathbf{x}}_i)), \mathbf{z}_i = \tilde{\mathbf{z}}_i / \|\tilde{\mathbf{z}}_i\|_2$ 7 // update class-prototypes  $\boldsymbol{\mu}_c := \text{Normalize}(\alpha \boldsymbol{\mu}_c + (1 - \alpha) \mathbf{z}_i), \ \forall c \in \{1, 2, \dots, C\}$ 8 // calculate compactness loss  $\mathcal{L}_{\text{comp}} = -\sum_{i=1}^{b} \log \frac{\exp(\mathbf{z}_{i}^{\top} \boldsymbol{\mu}_{c(i)} / \tau)}{\sum_{i=1}^{C} \exp(\mathbf{z}_{i}^{\top} \boldsymbol{\mu}_{j} / \tau)}$ 9 // calculate dispersion loss  $\mathcal{L}_{\text{dis}} = \frac{1}{C} \sum_{i=1}^{C} \log \frac{1}{C-1} \sum_{j=1}^{C} \mathbb{1}\{j \neq i\} e^{\mu_i^{\top} \mu_j / \tau}$ 10// calculate overall loss  $\mathcal{L} = \mathcal{L}_{dis} + \lambda_c \mathcal{L}_{comp}$ 11 // update the network weights update the weights in the encoder and the projection head 12

# **Experiment Setup**



#### **In-distribution datasets**



1) **CIFAR-10 [Krizhevsky 2012]** A dataset with 60,000 images composed of *five tasks* from *ten animal and vehicle classes*.



2) **CIFAR-100** [Krizhevsky 2012] A dataset with 60,000 images composed of 20 tasks from 100 generic object classes.

#### **OOD** test datasets



3) SVHN[Netzer 2011], Places365[Zhou 2017)], Textures[Cimpoi 2014], LSUN [Yu 2015], and iSUN [Xu 2015].

#### Test Time OOD Score

- **1)** Maha score (based on the Mahalanobis distance)
- 2) KNN score <sup>[2]</sup> (based on the Euclidean distance to the K-th nearest neighbor)



#### **Evaluation metrics.**

- 1) FPR95(the false positive rate of OOD samples when the true positive rate of ID samples is at 95%)
- 2) AUROC (the area under the receiver operating characteristic curve)
- **3) ID ACC (ID classification accuracy)**

[1] Lee et al., A Simple Unified Framework for Detecting OOD Saamples and Adversarial Attacks, NeurIPS 2018
 [2] Sun et al., Out-of-Distribution Detection with Deep Nearest Neighbbors, ICML 2022



#### **CIDER on Small-scale Datasets**

➤ CIFAR-100 (ID): methods with contrastive losses are trained for 500 epoch

➢ Backbone: ResNet-34

ID classification accuracy on CIFAR-100 (%)

ParNeC 模式识别与神经计算研究组 PAttern Recognition and NEural Computing

Method	SVHN		Places365		OOD Dataset LSUN		iSUN		Texture		Average		Method	ID ACC
	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	FPR↓	AUROC↑	w.o. contrast	ive loss
Without Contrastive Learning											MSP	74.59		
MSP	78.89	79.80	84.38	74.21	83.47	75.28	84.61	74.51	86.51	72.53	83.12	75.27	ODIN	74.59
ODIN	70.16	84.88	82.16	75.19	76.36	80.10	79.54	79.16	85.28	75.23	78.70	79.11	GODIN	74.92
Mahalanobis	87.09	80.62	84.63	73.89	84.15	79.43	83.18	78.83	61.72	84.87	80.15	79.53	Enorgy	74.50
Energy	66.91	85.25	81.41	76.37	59.77	86.69	66.52	84.49	79.01	79.96	70.72	82.55	Energy	74.39
GODIN	74.64	84.03	89.13	68.96	93.33	67.22	94.25	65.26	86.52	69.39	87.57	70.97	Mahalanobis	74.59
LogitNorm	59.60	90.74	80.25	78.58	81.07	82.99	84.19	80.77	86.64	75.60	78.35	81.74	w. contrasti	ve loss
					With C	ontrastive L	earning						CE + SimCLR	73.54
ProxyAnchor	87.21	82.43	70.10	79.84	37.19	91.68	70.01	84.96	65.64	84.99	66.03	84.78	SSD+	75.11
CE + SimCLR	24.82	94.45	86.63	71.48	56.40	89.00	66.52	83.82	63.74	82.01	59.62	84.15		75.11
CSI	44.53	92.65	79.08	76.27	75.58	83.78	76.62	84.98	61.61	86.47	67.48	84.83	ProxyAnchor	74.21
SSD+	31.19	94.19	77.74	79.90	79.39	85.18	80.85	84.08	66.63	86.18	67.16	85.90	KNN+	75.11
KNN+	39.23	92.78	80.74	77.58	48.99	89.30	74.99	82.69	57.15	88.35	60.22	86.14	CIDEP	75.25
CIDER	23.09	95.16	79.63	73.43	16.16	96.33	71.68	82.98	43.87	90.42	46.89	87.67	CIDEK	15.55

# Experiment

ParNeC 模式识别与神经计算研究组 PAttern Recognition and NEural Computing

#### **CIDER on Large-scale Datasets**

➤ ImageNet-100 (ID): finetune for 10 epoch



# Experiment



#### **CIDER** learns distinguishable representations.

Visualization of learned features by UMAP On CIFAR-10 (ID)



Figure 3: (a): UMAP (McInnes et al., 2018) visualization of the features when the model is trained with CE vs. CIDER for CIFAR-10 (ID). (b): CIDER makes OOD samples more separable from ID compared to CE (*c.f.* Table 4).



#### **CIDER** improves inter-class dispersion and intra-class compactness.

Compactness $(\mathcal{D}_{tr}^{in}, \boldsymbol{\mu}) = \frac{1}{C} \sum_{j=1}^{C} \frac{1}{n} \sum_{i=1}^{n} z_i^{\top} \boldsymbol{\mu}_j \mathbb{1}\{y_i = j\},\$ 

Dispersion $(\boldsymbol{\mu}) = \frac{1}{C} \sum_{i=1}^{C} \frac{1}{C-1} \sum_{j=1}^{C} \boldsymbol{\mu}_{i}^{\top} \boldsymbol{\mu}_{j} \mathbb{1}\{j \neq i\}.$ 

$$\uparrow \text{ Separability} = \frac{1}{|\mathcal{D}_{\text{test}}^{\text{ood}}|} \sum_{x \in \mathcal{D}_{\text{test}}^{\text{ood}}} \max_{j \in [C]} z_x^\top \boldsymbol{\mu}_j - \frac{1}{|\mathcal{D}_{\text{test}}^{\text{in}}|} \sum_{x' \in \mathcal{D}_{\text{test}}^{\text{in}}} \max_{j \in [C]} z_{x'}^\top \boldsymbol{\mu}_j,$$

Table 4: Compactness and dispersion of CIFAR-10 feature embedding, along with the separability *w.r.t.* each OOD test set. We convert cosine similarity to angular degrees for better readability.

Training Loss	Dispersion (ID) $\uparrow$	Compactness (ID)↓	<b>ID-OOD Separability</b> ↑ (in degree)					
	(in degree)	(in degree)	CIFAR-100	LSUN	iSUN	Texture	SVHN	AVG
Cross-Entropy SSD+ (SupCon loss)	67.17 75.50	24.53 22.08	7.11	14.57 28.55	13.70 25.70	13.76 33.45	$11.08 \\ 37.70$	12.04 29.86
CIDER (ours)	87.53	21.35	31.41	<b>48.37</b>	41.54	<b>39.60</b>	51.65	42.51

#### Experiment



#### **Ablation study**

- ➤ CIFAR-100 (ID): methods with contrastive losses are trained for 500 epoch
- ➢ Backbone: ResNet-34

Table 3: Ablation study on loss component. Results (in AUROC) are based on CIFAR-100 trained with ResNet-34. Training with only  $\mathcal{L}_{comp}$  suffices for ID classification. Inter-class dispersion induced by  $\mathcal{L}_{dis}$  is key to OOD detection.

Loss Co	mponents		ID ACC↑						
$\mathcal{L}_{ ext{comp}}$	$\mathcal{L}_{ ext{dis}}$	Places365	LSUN	iSUN	Texture	SVHN	AVG	ID Nee1	
$\checkmark$		79.63	85.75	84.45	87.21	91.33	85.67	75.19	
	$\checkmark$	54.76	69.81	54.99	44.26	46.48	54.06	2.03	
$\checkmark$	$\checkmark$	73.43	96.33	82.98	90.42	95.16	87.67	75.35	

# Thanks