
Fixed Non-negative Orthogonal Classifier: Inducing Zero-mean Neural Collapse with Feature Dimension Separation

Hoyong Kim, Kangil Kim
Artificial Intelligence Graduate School
Gwangju Institute of Science and Technology

■ Preliminaries

- **Neural Collapse (NC):** A recently discovered phenomenon that at the terminal phase of training, the **last-layer features of the same class** will **collapse** into **a single vertex**, and the vertices of all classes will be aligned with their **classifier prototypes** and be formed as **a simplex equi-angular tight frame (ETF)**
- **Layer-Peeled Model (LPM):** concentrate exclusively on the last-layer features \mathbf{h} and class weight vectors \mathbf{w} disregarding the encoder

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}_{ce}(\mathbf{h}_{k,i}, \mathbf{w}_k),$$

$$s. t. \|\mathbf{w}_k\|^2 \leq E_W, \forall 1 \leq k \leq K,$$

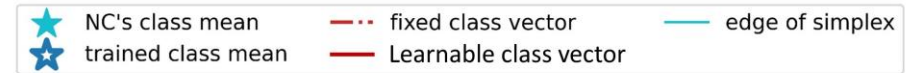
$$\|\mathbf{h}_{k,i}\|^2 \leq E_H, \forall 1 \leq k \leq K, \forall 1 \leq i \leq n_k$$

The LPM has been proven to achieve global optimality when satisfying NC properties:

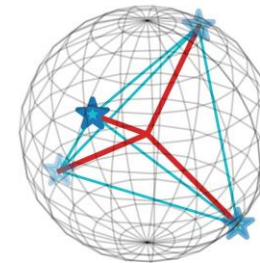
$$\mathbf{h}_{k,i}^* = C\mathbf{w}_k^* = C'\mathbf{m}_k^*,$$

where the matrix $[\mathbf{m}_1^*, \dots, \mathbf{m}_K^*]$ forms a K -simplex ETF

■ Motivation



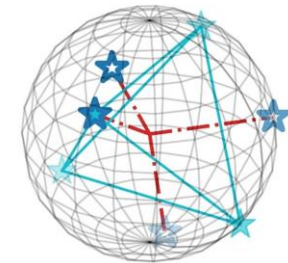
Learnable Classifier



Neural collapse **occurs**

class vectors -> simplex ETF
trained class means -> simplex ETF

Fixed Classifier (not a simplex)



Collapse **occurs in a different shape**

class vectors -> non-simplex ETF
trained class means -> non-simplex ETF

How does the **collapse** between class means and class weight vectors occur in the fixed classifier when their shape is **not a simplex ETF**?

Research Problem

- When the LPM with a fixed classifier which shape is non-simplex ETF, NC cannot explain the collapse phenomenon of it.

Solution

Non-negativity and Orthogonality

: These two additional constraints in the LPM, it can achieve the global optimality even in inducing the max-margin decision

→ **Fixed Non-negative Orthogonal Classifier**

Zero-mean Neural Collapse (ZNC)

: when the LPM with a fixed non-negativity orthogonal classifier achieves the global optimality, it satisfies a different collapse properties.

→ To explain it, we propose a zero-mean neural collapse

Orthogonal Layer-Peeled Model (OLPM): the layer-peeled model with a fixed non-negative orthogonal classifier \mathbf{Q}^*

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathcal{L}_{ce}(\mathbf{h}_{k,i}, \mathbf{Q}^*),$$

$$s. t. \|\mathbf{h}_{k,i}\|^2 \leq E_H, \text{ and } \sum_{j \neq k} \mathbf{h}_{k,i}^\top \mathbf{q}_j^* \geq 0,$$

$$\forall 1 \leq k \leq K, \forall 1 \leq i \leq n_k$$

The OLPM achieves global optimality when satisfying ZNC properties:

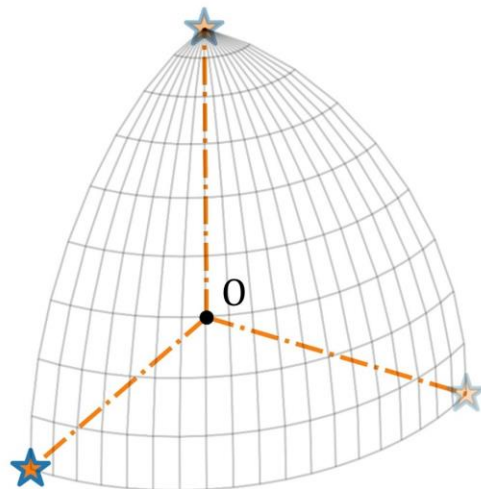
$$\mathbf{h}_{k,i}^* \mathbf{q}_{k'}^* = (K - 1) \delta_{k,k'}$$

where $\delta_{k,k'}$ is a kronecker delta function

Definition 1 (Non-negative Orthogonal Classifier). A *non-negative orthogonal classifier* has a partial orthogonal weight matrix $\mathbf{Q} \in \mathbb{R}_{\geq 0}^{D \times K}$, which satisfies below properties:

$$\mathbf{Q}^T \mathbf{Q} = \mathbf{I}_K, \quad \text{s.t. } Q_{i,j} \geq 0, \forall 1 \leq i \leq D, \forall 1 \leq j \leq K,$$

where $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ is the identity matrix and $Q_{i,j}$ is the (i, j) element of \mathbf{Q}



* The properties of Zero-mean Neural Collapse

: The only difference with NC is that **class means are centered to the origin, not their global mean**, as still satisfying the properties of neural collapse

- FNO classifier invokes *feature dimension separation* (FDS), which **reduces the interference between class weight vectors**

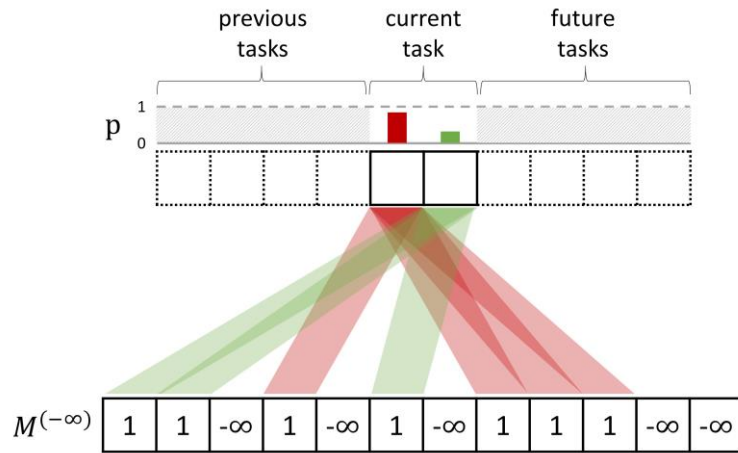
→ **FDS is useful in continual learning and imbalanced learning !**

Definition 2 (Feature Dimension Separation). Let $\mathbf{q}_k = \{q_j\}_{1 \leq j \leq D}$ the k -th class weight vector in the fixed non-negative orthogonal classifier and $\mathbb{J}_k = \{j \mid q_j > 0, 1 \leq j \leq D\}$ an index set of \mathbf{q}_k where q_j is not zero. Then, as the definition of FNO classifier, any index set of class weight vectors has disjoint to any other class weight vector and we call this phenomenon as feature dimension separation, i.e.,

$$\mathbb{J}_k \cap \mathbb{J}_{k'} = \emptyset, \quad \forall k \neq k',$$

which means that *the features' elements used for deciding the confidence to any class lose their utility to any other classes.*

$$\mathbf{p} = \text{WeightedSoftmax} \left(M^{(-\infty)} \odot \text{MatMul}(\mathbf{Q}, \mathbf{H}) \right)$$



Definition 3 (Masked Softmax). To remove the class-wise interference of specific classes in softmax, the masked softmax multiplies a negative infinity mask $\mathbf{M}^{(-\infty)}$ to output vectors. When getting rid of k -th class's interference from i -th input sample, k -th element of the output vector is multiplied by negative infinity values, i.e.,

$$\mathbf{M}_i^{(-\infty)} = (m_j)_{1 \leq j \leq K}$$

$$\mathbf{p} = \text{Softmax} \left(\mathbf{M}_i^{(-\infty)} \odot (\mathbf{W}^\top \mathbf{h} + \mathbf{b}) \right),$$

where $m_j = -\infty$ if $j = k$ otherwise 1 and \odot means the Hadamard product.

Algorithm 1 Masked and Weighted Softmax with FNO classifier in a Task T_t

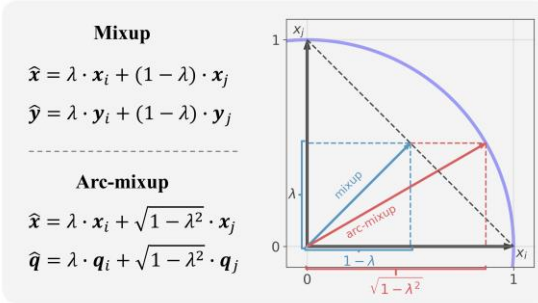
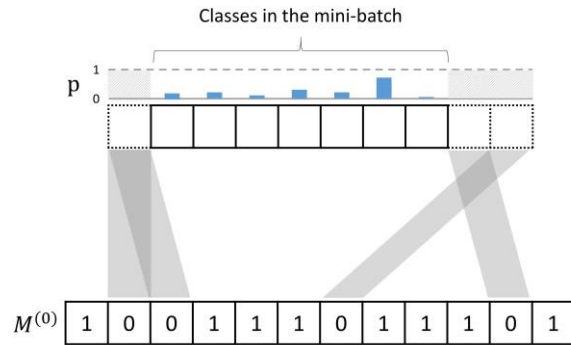
Require: $(\mathbf{X}_t, \mathbf{Y}_t), \mathbf{Q}$

Ensure: $\mathbf{P} \in \mathbb{R}^{N_t \times K}$

- 1: $\mathbf{H} \leftarrow \text{RELU}(f_\theta(\mathbf{X}_t))$ ▷ get features from \mathbf{X}_t as Eq. 1
 - 2: $\mathbb{K} = \{c_i \mid c_i \text{ of } \mathbf{X}_t\}$ ▷ initialize a set of class labels in the task T_t .
 - 3: $M^{(-\infty)} = (m_{i,j})_{1 \leq i \leq N_t, 1 \leq j \leq K}$, where $m_{i,j} = -\infty$ if $j \notin \mathbb{K}$ otherwise 1 ▷ initialize $M^{(-\infty)}$
 - 4: $\mathbf{P} = \text{W-SOFTMAX}(M^{(-\infty)} \odot \text{MATMUL}(\mathbf{Q}, \mathbf{H}))$ ▷ get the confidence of \mathbf{H}
-

$$\mathbf{p} = \text{MatMul}(\hat{\mathbf{Q}}, \text{LayerNorm}(M^{(0)} \odot \hat{\mathbf{H}}))$$

\mathbf{H} : last-layer features
 \mathbf{Q} : non-negative and orthogonal matrix
 $M^{(-\infty)}$: negative infinite mask $M^{(0)}$: zero mask



Definition 4 (Arc-mixup). Arc-mixup is an interpolation-based method when all last-layer features and class weight vectors are located on the same hypersphere, while keeping the scale of mixed class weight vector $\hat{\mathbf{q}}$, i.e.,

$$\hat{\mathbf{x}} = \lambda \cdot \mathbf{x}_i + \sqrt{1 - \lambda^2} \cdot \mathbf{x}_j$$

$$\hat{\mathbf{q}} = \lambda \cdot \mathbf{q}_i + \sqrt{1 - \lambda^2} \cdot \mathbf{q}_j$$

$$\mathcal{L}_{cls}(\hat{\mathbf{x}}, \hat{\mathbf{q}}) = -\log \hat{\mathbf{q}}^T \hat{\mathbf{h}}$$

and $\hat{\mathbf{h}}$ is the last-layer feature of $\hat{\mathbf{x}}$. This means that $\hat{\mathbf{q}}$ is still located on the hypersphere

Algorithm 2 Arc-mixup with FNO classifier and feature masking in a mini-batch \mathbb{B}

Require: $(\mathbf{X}, \mathbf{Y}) \in \mathbb{B}, \mathbf{Q}$

Ensure: $\mathbf{P} \in \mathbb{R}^{|\mathbb{B}| \times K}$

1: $(\hat{\mathbf{X}}, \hat{\mathbf{Q}}) \leftarrow \text{ArcMixup}(\mathbf{X}, \mathbf{Q})$

2: $\hat{\mathbf{H}} \leftarrow \text{RELU}(f_\theta(\hat{\mathbf{X}}))$

3: $\mathbb{K} = \{c_i \mid c_i \in \mathbb{B}, \forall 1 \leq i \leq |\mathbb{B}|\}$

4: $\hat{\mathbb{J}} = \bigcup_{k \in \mathbb{K}} \mathbb{J}_k$

5: $M^{(0)} = (m_{i,j})_{1 \leq i \leq |\mathbb{B}|, 1 \leq j \leq D}$, where $m_{i,j} = \mathbf{1}_{j \in \hat{\mathbb{J}}}$

6: $\mathbf{P} = \text{MATMUL}(\hat{\mathbf{Q}}, \text{LAYERNORM}(M^{(0)} \odot \hat{\mathbf{H}}))$

▷ mixup input samples and class vectors as Eq. 11

▷ get features from $\hat{\mathbf{X}}$ as Eq. 1

▷ initialize a set of class labels in the mini-batch

▷ initialize an index set including all index sets of \mathbf{Q}

▷ initialize a zero mask $M^{(0)}$

▷ get the confidence of $\hat{\mathbf{H}}$

Experimental Results

In classification results for standard CL bench-marks, our method demonstrated superior performance in all class-incremental learning settings (*Max: +5.23 in S-TinyImageNet with buffer size 200*)

\mathcal{B}	Method			S-MNIST		S-CIFAR-10		S-CIFAR-100		S-Tiny-ImageNet	
	RM	Cif	M	Class-IL	Task-IL	Class-IL	Task-IL	Class-IL	Task-IL	Class-IL	Task-IL
200	ER	FC	✓	82.98 _{1.03}	98.13 _{0.16}	61.75 _{6.07}	91.39 _{2.13}	28.51 _{0.44}	68.51 _{0.87}	15.47 _{0.67}	44.11 _{0.50}
	ER	FNO	✓	84.26 _{1.16}	98.45 _{0.19}	63.84 _{1.47}	92.03 _{0.52}	32.43 _{0.58}	71.34 _{1.17}	17.31 _{0.74}	44.76 _{0.90}
				+1.28	+0.32	+2.09	+0.64	+3.92	+2.83	+1.84	+0.65
	DER++	FC	✓	84.45 _{0.88}	99.03 _{0.09}	66.35 _{1.52}	93.17 _{0.54}	28.57 _{1.11}	74.02 _{0.76}	13.21 _{0.56}	49.75 _{0.99}
	DER++	FNO	✓	86.27 _{0.88}	99.11 _{0.08}	67.53 _{1.25}	93.98 _{0.39}	30.70 _{1.16}	74.11 _{0.96}	18.44 _{0.94}	53.06 _{0.67}
				+1.82	+0.08	+1.18	+0.81	+2.13	+0.09	+5.23	+3.31
500	ER	FC	✓	89.35 _{0.59}	99.20 _{0.16}	70.64 _{1.28}	94.22 _{0.41}	35.68 _{0.89}	74.77 _{0.71}	20.43 _{0.38}	53.21 _{0.84}
	ER	FNO	✓	89.42 _{0.72}	99.16 _{0.17}	71.43 _{0.95}	94.38 _{0.43}	39.80 _{0.68}	76.52 _{0.86}	22.41 _{0.57}	52.60 _{0.58}
				+0.07	-0.04	+0.79	+0.16	+4.12	+1.75	+1.98	-0.61
	DER++	FC	✓	83.10 _{1.22}	99.08 _{0.09}	71.85 _{3.76}	94.28 _{1.49}	37.80 _{0.92}	80.52 _{0.60}	17.71 _{0.58}	59.86 _{1.08}
	DER++	FNO	✓	86.75 _{0.75}	99.00 _{0.10}	74.77 _{0.66}	95.56 _{0.16}	40.81 _{0.70}	80.61 _{0.46}	22.45 _{0.36}	59.87 _{1.91}
				+3.65	-0.08	+2.92	+1.28	+3.01	+0.09	+4.74	+0.01
5120	ER	FC	✓	93.51 _{0.60}	99.38 _{0.12}	82.63 _{1.34}	96.45 _{0.27}	52.95 _{0.73}	84.20 _{0.58}	35.73 _{0.41}	67.50 _{0.53}
	ER	FNO	✓	93.98 _{0.39}	99.47 _{0.09}	82.88 _{1.35}	96.79 _{0.38}	57.02 _{0.52}	85.66 _{0.42}	36.90 _{0.41}	66.86 _{0.32}
				+0.47	+0.09	+0.25	+0.34	+4.07	+1.46	+1.17	-0.64
	DER++	FC	✓	93.75 _{0.23}	99.62 _{0.05}	84.71 _{0.65}	96.78 _{0.16}	58.18 _{0.43}	87.97 _{0.33}	34.72 _{0.46}	72.40 _{0.25}
	DER++	FNO	✓	94.26 _{0.24}	99.59 _{0.05}	85.65 _{0.38}	97.20 _{0.13}	58.82 _{0.43}	87.35 _{0.45}	38.95 _{0.71}	72.70 _{0.27}
				+0.51	-0.03	+0.94	+0.42	+0.64	-0.62	+4.23	+0.30

Experimental Results

In imbalanced learning on CIFAR10/100-LT, ImageNet-LT and Places-LT, our method performed better than other methods (*Max: +9.90 in Places-LT*)

Method			Reference	CIFAR10-LT			CIFAR100-LT		
Aug	Clf	\mathcal{L}		100	50	10	100	50	10
mixup	FC	CE	(Yang et al., 2022b)	73.90 _{0.30}	79.30 _{0.20}	87.80 _{0.10}	43.00	48.10	-
B-mixup	FC	CE	(Zhang et al., 2022b)	78.70	-	89.60	-	-	-
mixup	ETF	CE	(Yang et al., 2022b)	67.00 _{0.40}	77.20 _{0.30}	87.00 _{0.20}	-	-	-
mixup	ETF	DR	(Yang et al., 2022b)	76.50 _{0.30}	81.00 _{0.20}	87.70 _{0.20}	45.30	50.40	-
mixup	FC	CE	(reproduced.) [†]	74.24 _{0.44}	80.00 _{0.54}	89.08 _{0.32}	43.80 _{0.42}	49.57 _{0.37}	63.90 _{0.33}
mixup	ETF	DR	(reproduced.) [†]	75.18 _{0.49}	80.17 _{0.31}	87.29 _{0.23}	45.45 _{0.38}	50.67 _{0.37}	62.84 _{0.38}
arc-mixup	FNO	CE	(reproduced.) [†]	82.59 _{0.26}	85.13 _{0.25}	89.50 _{0.14}	49.26 _{2.82}	54.44 _{2.32}	63.14 _{3.82}

Method			Reference	ImageNet-LT (ResNet50)				Places-LT	
Aug	Clf	\mathcal{L}		Many	Median	Few	All	ResNet152	ResNet152 (FT)
mixup	FC	CE	(Yang et al., 2022b)*	-	-	-	44.30	-	-
mixup	ETF	DR	(Yang et al., 2022b)*	-	-	-	44.70	-	-
-	FC	CE	(reproduced) [†]	66.53 _{0.18}	40.34 _{0.45}	12.03 _{0.30}	45.88 _{0.31}	22.62 _{0.29}	24.16 _{0.41}
mixup	FC	CE	(reproduced) [†]	67.40 _{0.55}	38.74 _{0.83}	9.12 _{0.40}	45.03 _{0.64}	22.10 _{0.27}	24.83 _{1.19}
mixup	ETF	DR	(reproduced.) [†]	64.17 _{0.27}	22.12 _{0.38}	0.57 _{0.16}	34.96 _{0.27}	23.11 _{0.16}	25.51 _{0.08}
arc-mixup	FNO	CE	(reproduced) [†]	59.46 _{0.51}	43.54 _{0.26}	24.48 _{0.55}	46.60 _{0.40}	30.07 _{0.40}	34.06 _{0.10}

▪ Contributions

- **Zero-mean Neural Collapse**
: we propose a *zero-mean neural collapse* to analyze the collapse phenomenon in training classification model with the fixed classifier
- **Fixed Non-negative Orthogonal Classifier**
: we propose a *fixed non-negative orthogonal classifier* and prove its theoretical benefits in orthogonal layer-peeled model with the zero-mean neural collapse
- **Benefits of Feature Dimension Separation**
: we demonstrate the impacts of the proposed methods with masked softmax in continual learning and arc-mixup in imbalanced learning

Thank You

<https://github.com/GIST-IRR/FNO-classifier>

hoyong.kim.21@gm.gist.ac.kr