

Neural Collapse의 디자인, CE와 MSE의 차이 분석

MSE의 차이 새로운 decomposition.

MSE classifier의 deviation capture

### \* Introduction.

MLP, SGD, BN, CE 등은 관행적으로 쓰이고 이에 대한 분석이 중요하다.

논문에서는 classifier의 마지막 layer의 대해 이를적으로 조사. (C)는 class, N은 sample)

classifier의 input(아키텍처의 feature  $h \in \mathbb{R}^p$ )

$$\hookrightarrow w_c \in \mathbb{R}^p, b_c \in \mathbb{R}$$

이미 디蹲 error는  $\text{Error} = \text{Ave}_{i,c} \mathbf{1}\{c \neq \arg \max_{c'} (\langle w_{c'}, h_{i,c} \rangle + b_{c'})\}$ ,

CE는 이미 디蹲  $\text{CE} = -\text{Ave}_{i,c} \log \frac{\exp\{\langle w_c, h_{i,c} \rangle + b_c\}}{\sum_{c'=1}^C \exp\{\langle w_{c'}, h_{i,c} \rangle + b_{c'}\}}$ ,

이전 논문에서 overparameterized된 classifier는 test의 디蹲 손상시키지 않고 train을 memorize

memorize 이후 계속 훈련하면 성능 향상을 할 수 있음.

이미 대해 Error 0을 갖는 zero-CE-loss를 향해 가는 것을 Terminal Phase of Training (TPT)

TPT 동안 Neural Collapse (NC) 나타남.

### - Neural Collapse

Feature Global Mean :  $\mu_G = \text{Ave}_{i,c} h_{i,c}$

Feature Class Mean :  $\mu_c = \text{Ave}_{i,c} h_{i,c}, c = 1, \dots, C$

Feature within-Class Covariance :  $\Sigma_w = \text{Ave}_{i,c} (h_{i,c} - \mu_c)(h_{i,c} - \mu_c)^T$

Feature between-Class Covariance :  $\Sigma_B = \text{Ave}_i (\mu_c - \mu_G)(\mu_c - \mu_G)^T$

훈련 epoch 틱이 지나면收敛하는 (global limiting behavior)가 있다.

#### NC1. Within-class variability collapse

$$\Sigma_B^\dagger \Sigma_W \rightarrow \mathbf{0}, \quad \text{t.e. Moore-Penrose pseudoinverse}$$

#### NC2. Convergence to Simplex ETF

$$\frac{\langle \mu_c - \mu_G, \mu_{c'} - \mu_G \rangle}{\|\mu_c - \mu_G\|_2 \|\mu_{c'} - \mu_G\|_2} \rightarrow \begin{cases} 1, & c = c' \\ \frac{-1}{C-1}, & c \neq c' \end{cases}$$

$$\|\mu_c - \mu_G\|_2 - \|\mu_{c'} - \mu_G\|_2 \rightarrow 0 \quad \forall c \neq c'$$

#### NC3. Convergence to self-duality

$$\frac{w_c}{\|w_c\|_2} - \frac{\mu_c - \mu_G}{\|\mu_c - \mu_G\|_2} \rightarrow 0$$

#### NC4. Simplification to nearest class center

$$\arg \max_{c'} \langle w_{c'}, h \rangle + b_{c'} \rightarrow \arg \min_{c'} \|h - \mu_{c'}\|_2$$

NC2는 ETF라는 단순한 기하학적 구조로 수렴을 capture

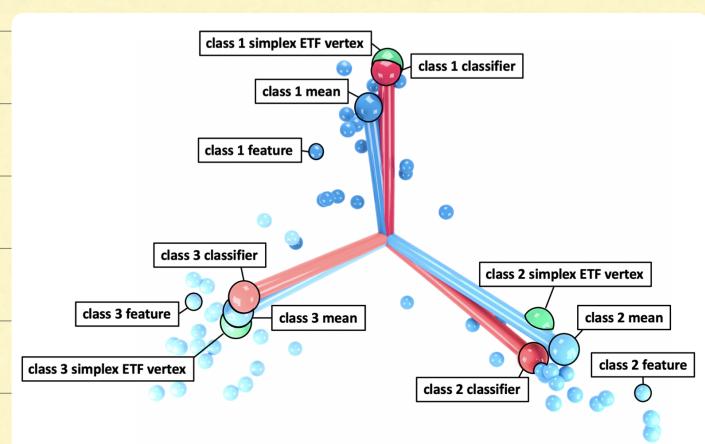
ETF는 길이가 같고 최대로 분리된 pair-wise angle을 갖는 vector  $\{v_c\}_{c=1}^C$ 의 collection이다.

$P > C$ 인 setting에서 maximal angle은

$$\frac{\langle v_c, v_{c'} \rangle}{\|v_c\|_2 \|v_{c'}\|_2} = \begin{cases} 1, & \text{for } c = c' \\ -\frac{1}{C-1}, & \text{for } c \neq c' \end{cases}$$

여기서 ETF는 simplex ETF로 불림.

$C$ 가 증가하면 orthogonal이 되도록됨.



### - DeepNet classification with MSE Loss

본류 net은 보통 CE로 훈련되지만 몇몇은 MSE로 훈련함.

$$\begin{aligned}\mathcal{L}(\mathbf{W}, \mathbf{b}, \mathbf{H}) &= \frac{1}{2} \operatorname{Ave}_{i,c} \|\mathbf{W}\mathbf{h}_{i,c} + \mathbf{b} - \mathbf{y}_{i,c}\|_2^2 + \frac{\lambda}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{b}\|_2^2) \\ &= \frac{1}{2CN} \|\mathbf{WH} + \mathbf{b}\mathbf{1}_{CN}^\top - \mathbf{Y}\|_F^2 + \frac{\lambda}{2} (\|\mathbf{W}\|_F^2 + \|\mathbf{b}\|_2^2),\end{aligned}$$

$$\mathbf{H} \in \mathbb{R}^{P \times CN}, \mathbf{Y} \in \mathbb{R}^{C \times CN}, \mathbf{W} \in \mathbb{R}^{C \times P}, \mathbf{b} \in \mathbb{R}^C$$

### - Contribution.

1. MSE의 대체적 세로운 decomposition.  $\mathcal{L} = \mathcal{L}_{NC1} + \mathcal{L}_{NC2/3} + \mathcal{L}_{LS}^\perp$

↳ classifier  $\mathbf{W}$ 는 MSE의 다른 optimal classifier  $\mathbf{W}_{LS}$ 와 같음

2. real dataset의 decomposition의 경향적 특성과  $\mathcal{L}_{LS}^\perp$ 가 훈련 중에 부시될 수 있음을 보임

↳ central path:  $\mathcal{L}_{LS}^\perp = \text{zero}$

3. central path의 대수 invariant property임을 밝힘

↳  $\mathbf{X} = \sum_w^{-\frac{1}{2}} \mathbf{H}$  (renormalized features)

4. central의 대수 renormalized feature의 gradient flow를 연구하고

NC를 의인화하는 closed-form dynamics 도출

dynamics는 initial renormalized feature class-mean의 singular value decomposition에 의해 명시적

## \* Decomposition of MSE Loss

MSE의 경우 NC 분석을 위한 decomposition

$W$ 와  $b$ 를 합친다면,  $\tilde{W} = [W, b] \in \mathbb{R}^{C \times (P+1)}$ ,  $\tilde{h}_{i,c} = [h_{i,c}; 1] \in \mathbb{R}^{P+1}$

$$\mathcal{L}(\tilde{W}, \tilde{H}) = \frac{1}{2} \operatorname{Ave}_{i,c} \|\tilde{W}\tilde{h}_{i,c} - \mathbf{y}_{i,c}\|_2^2 + \frac{\lambda}{2} \|\tilde{W}\|_F^2.$$

$\tilde{h}$ 로 확장된 NC를 보면,

$$\tilde{\Sigma}_T = \operatorname{Ave}_{i,c} (\tilde{h}_{i,c} - \tilde{\mu}_G)(\tilde{h}_{i,c} - \tilde{\mu}_G)^\top \in \mathbb{R}^{(P+1) \times (P+1)}$$

$$\tilde{M} = [\tilde{\mu}_1, \dots, \tilde{\mu}_C] \in \mathbb{R}^{(P+1) \times C}.$$

Proposition 1.

고정된  $\tilde{H}$ 에서 MSE의 case optimal classifier는

$$\tilde{W}_{LS} = \frac{1}{C} \tilde{M}^\top (\tilde{\Sigma}_T + \tilde{\mu}_G \tilde{\mu}_G^\top + \lambda I)^{-1},$$

$\hookrightarrow \tilde{H}$ 이면 dependent.

Theorem 1. - Appendix B

MSE Loss  $\mathcal{L}(\tilde{W}, \tilde{H}) = \mathcal{L}_{LS}(\tilde{H}) + \mathcal{L}_{LS}^\perp(\tilde{W}, \tilde{H})$ 로 분해 가능.

$$\mathcal{L}_{LS}(\tilde{H}) = \frac{1}{2} \operatorname{Ave}_{i,c} \|\tilde{W}_{LS} \tilde{h}_{i,c} - \mathbf{y}_{i,c}\|_2^2 + \frac{\lambda}{2} \|\tilde{W}_{LS}\|_F^2,$$

$\tilde{W}$ 의 independent

$$\mathcal{L}_{LS}^\perp(\tilde{W}, \tilde{H}) = \frac{1}{2} \operatorname{tr} \left\{ (\tilde{W} - \tilde{W}_{LS}) \left( \tilde{\Sigma}_T + \tilde{\mu}_G \tilde{\mu}_G^\top + \lambda I \right) (\tilde{W} - \tilde{W}_{LS})^\top \right\}.$$

직관적으로  $\tilde{W}_{LS}$ 는  $\tilde{W}$  보면  $\tilde{H}$ 의 연관이 있다.

$\mathcal{L}_{LS}^\perp$ 은 non-negative I.  $\tilde{W} = \tilde{W}_{LS}$  일 때면 0이다.

$\therefore \mathcal{L}_{LS}^\perp$ 은  $\tilde{W}$ 와  $\tilde{W}_{LS}$  간의 distance로 볼 수 있음

$\mathcal{L}_{LS}^\perp$ 은 network의 활동을 capture

$\mathcal{L}_{LS}$ 를 추가적으로 decompose.  $\rightarrow NC1, NC2$ 을 보기 위해.

### Theorem 2 - Appendix C.

Theorem 1의  $\mathcal{L}_{LS}(\tilde{\mathbf{H}})$ 는 어떻게 decompose가능.

$$\mathcal{L}_{LS}(\tilde{\mathbf{H}}) = \mathcal{L}_{NC1}(\tilde{\mathbf{H}}) + \mathcal{L}_{NC2/3}(\tilde{\mathbf{H}})$$

$$\begin{aligned}\mathcal{L}_{NC1}(\tilde{\mathbf{H}}) &= \frac{1}{2} \text{tr} \left\{ \tilde{\mathbf{W}}_{LS} \left[ \tilde{\Sigma}_W + \lambda \mathbf{I} \right] \tilde{\mathbf{W}}_{LS}^\top \right\}, \\ \mathcal{L}_{NC2/3}(\tilde{\mathbf{H}}) &= \frac{1}{2C} \|\tilde{\mathbf{W}}_{LS} \tilde{\mathbf{M}} - \mathbf{I}\|_F^2.\end{aligned}$$

$\mathcal{L}_{NC2/3}$ 의 class-mean과 MSE-optimal classifier와 관련 있는 특수

Minimizing  $\mathcal{L}_{NC2/3}$ 은 class-mean과 classifier가 동일한 Simplex ETF matrix로 이동한다.

NC1은 mean과는 independent.

But then classifier가 potentially large ETF matrix와 연결되었지만  $\sum_w$ 를 0으로 유도해,  $\mathcal{L}_{NC1}$ 을 줄일 수 있다.

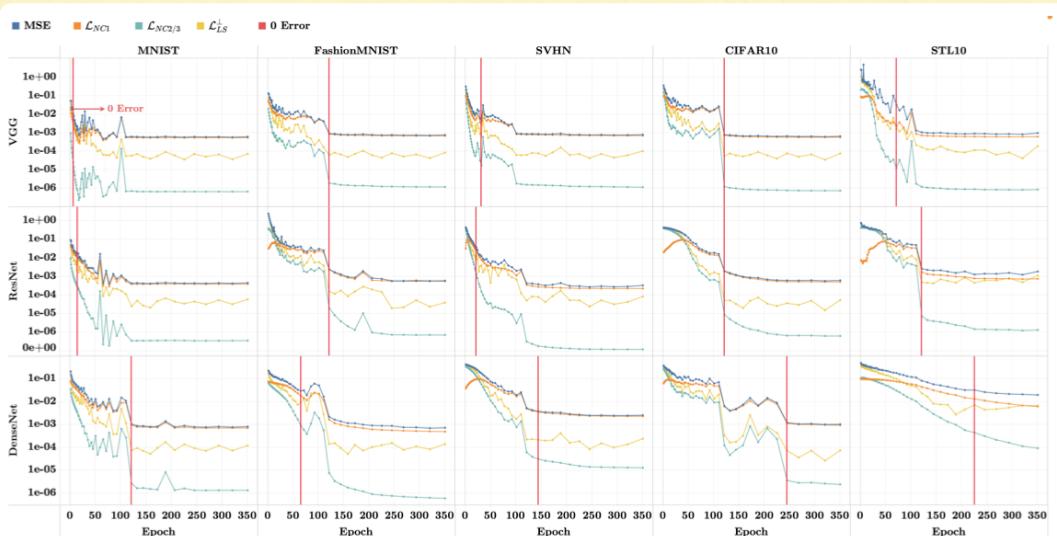


Figure 2: *Decomposition of MSE loss*: Each array column shows a benchmark image classification dataset while each row shows a canonical deep net architecture trained with MSE loss. The red vertical line indicates the epoch at which zero training error was achieved. In each array cell, we plot terms of the MSE loss decomposition  $\mathcal{L}(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) = \mathcal{L}_{NC1}(\tilde{\mathbf{H}}) + \mathcal{L}_{NC2/3}(\tilde{\mathbf{H}}) + \mathcal{L}_{LS}^\perp(\tilde{\mathbf{W}}, \tilde{\mathbf{H}})$  from Section 2. Starting from an early epoch in training,  $\mathcal{L}_{LS}^\perp(\tilde{\mathbf{W}}, \tilde{\mathbf{H}})$  becomes negligible compared to the dominant term,  $\mathcal{L}_{NC1}(\tilde{\mathbf{H}})$ , implying  $\mathcal{L}_{LS}^\perp(\tilde{\mathbf{W}}, \tilde{\mathbf{H}}) \ll \mathcal{L}_{LS}(\tilde{\mathbf{H}}) = \mathcal{L}_{NC1}(\tilde{\mathbf{H}}) + \mathcal{L}_{NC2/3}(\tilde{\mathbf{H}})$ , i.e. the features and classifiers are *effectively on the central path* during TPT. Note that  $\mathcal{L}_{NC2/3}(\tilde{\mathbf{H}})$  diminishes the fastest among all the terms: Intuitively, this shows that the network primarily focuses on distributing the feature class-means into a “uniform” Simplex ETF configuration (NC1)-(NC2) early on and, from there, compresses the activations towards their class-means, i.e. (NC1), as much as possible. Further experimental details are in Appendix A. Outlier behavior is discussed in Appendix A.7.

$\rightarrow \mathcal{L}_{LS}^\perp$ 는  $\mathcal{L}_{LS}$ 에 비해 무시할 만한 수준

$$\hookrightarrow \mathcal{L}(\tilde{W}, \tilde{H}) \approx \mathcal{L}_{LS}(\tilde{H})$$

↓ central path

$$\mathcal{P} = \left\{ (\tilde{W}_{LS}(\tilde{H}), \tilde{H}) \mid \tilde{H} \in \mathbb{R}^{(P+1) \times CN} \right\},$$

$\tilde{W}_{LS}$ 은  $\tilde{H}$ 에만 dependent 하기 때문에

$\tilde{W}$ 은 central path의 일환은  $\tilde{W}_{LS}$ 와 같다.

\* Exact Closed-Form Analysis on Central Path.