# Paper: The details matter: preventing class collapse in supcon

## Notations:

Labeled input data: $D = \{(x_i, y_i)\}_{i=1}^{N}$

$$(x, y) \in P$$

$$x \in \mathcal{X} \in \mathbb{R}^{\text{data dim.}}$$

$$y \in \mathcal{Y} = \{1, 2, \cdots K\}$$

Data label $\quad h(x) \in \mathcal{Y} \quad // \quad P(y \mid x)$

$$P(y = i) = \frac{1}{k}$$

Target : to learn a model $\hat{P}(y \mid x)$

$\underbrace{\text{Data point belongs to categories beyond labels}}$

<span style="color:red">STRATA</span>

Strata as latent labels $z \in \mathbb{Z} = \{1, 2, \cdots c\}$

$\mathbb{Z}$ divided into disjoint subset : $S_1, S_2 \cdots S_k$

$$z \in S_k, \quad y = k$$

$S(c)$ denotes deterministic label $c$.

1st strata is sampled $R(z)$

$x$ is sampled $P_z = P(\cdot \mid z)$

label is $y = S(z)$

# SupCon and collapse embedding

Similarity

$$\sigma(x, x') = f(x)^T f(x) / c$$

$B \to$ set of batches of labeled dataset on $D$

positive

$$\textcircled{P}(i, B) = \{ p \in B \setminus i : h(p) = h(i) \}$$

SupCon loss :

$$\hat{L}_{se}(f, x_i, B) = \frac{-1}{|P(i, B)|} \sum_{p \in P(i, B)} \log \frac{\exp(\sigma(x_i, x_p))}{\sum_{a \in B \setminus i} \exp(\sigma(x_i, x_a))}$$

positive pair.

Negative pair.

class collapse : simplex embedding Scenario.

if $h(x) = i$ then $f(x) = v_i \quad \forall x \in B$

$\{v_i\}_{i=1}^K \to$ forms regular simplex.

with

properties $\begin{cases} i) \sum_{i=1}^k v_i = 0 \\ \\ ii) \|v_i\|_2 = 1 \end{cases}$

(iii) $\exists\, c_k \in \mathbb{R}$ s.t. $\vec{v_i}^T \vec{y_j} = c_k$ for if $\vec{j}$

## End model

Linear Classifier: $W \in \mathbb{R}^{k \times d}$

$$\|w_y\|_2 \leq 1 \; ; \; y \in \mathcal{Y}$$

Empirical loss:

$$\hat{\mathcal{L}}(w, D) = \sum_{z_i \in D} -\log \frac{\exp\left(f(x_i)^T w_{h(x_i)}\right)}{\sum_{j=1}^{k} \exp\left(f(x_i)^T w_j\right)}$$

Prediction:

$$\hat{P}(y|x) = \hat{P}(y|f(x))$$

generalized error:

$$\mathcal{L}(x, y, f) = \mathbb{E}_{x,y}\left[-\log \hat{P}(y|f(x))\right]$$

# Methodologies

1. class collapse minimize $\mathcal{L}(x, z, f) \; \forall \; x$

losses
strata:
$\Big\{$ when , i) $p(y = h(x) \mid x) = 1$

ii) $p(z \mid y) = \frac{1}{m} \; ; \; z \in S_{h(x)}$

iii) $p(x \mid z) = p(x \mid y)$  // no strata distinction

## Modified Contrastive loss

$\mathcal{L}_{spread} = \alpha \, \mathcal{L}_{attract} + (1 - \alpha) \, \mathcal{L}_{repel}.$

Negative examples : $N(i, B) = \{ a \in B \setminus i : h(a) = h(i) \}$

$\hat{\mathcal{L}}_{att}(f, x_i, B) = \frac{-1}{|P(i, B)|} \sum_{p \in P(i, B)} \log \frac{\exp(\sigma(x_i, x_p))}{\exp(\sigma(x_i, x_p))}$

$+ \sum_{a \in N(i, B)} \exp(\sigma(x_i, x_a))$

$\hat{\mathcal{L}}_{rep} = - \log \frac{\exp(\sigma(x_i, x_i^{aug}))}{\sum_{p \in P(i, B)} \exp(\sigma(x_i, x_p))}$ $\begin{bmatrix} \text{typical} \\ \text{solution} \end{bmatrix}$

spread the positive around.