

Generalization Analysis for Supervised Contrastive Representation Learning under Non-IID Settings

Nong Minh Hieu, Antoine Ledent

School of Computing and Information Systems, SMU Singapore

Contrastive learning overview

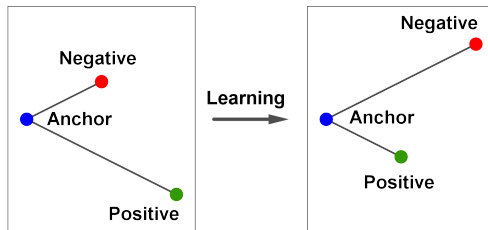


Figure: Visual illustration of contrastive learning

Contrastive Learning Overview

Problem Settings: Let \mathcal{X} be the input space and \mathcal{C} be the finite set of labels. Furthermore, denote

- ρ as the class distribution over \mathcal{C} .
- \mathcal{D}_c as the class-conditional distribution given label $c \in \mathcal{C}$.
- $\bar{\mathcal{D}}_c$ as the class-conditional distribution given labels $\mathcal{C} \setminus \{c\}$.

Definition (Unsupervised Risk): Given a contrastive loss $\ell : \mathbb{R}^k \rightarrow [0, \mathcal{M}]$, the unsupervised risk is defined as

$$L_{\text{un}}(f) = \mathbb{E}_{c \sim \rho} \mathbb{E}_{\substack{\mathbf{x}, \mathbf{x}^+ \sim \mathcal{D}_c^2 \\ \mathbf{x}_{1:k}^- \sim \bar{\mathcal{D}}_c^k}} \left[\ell \left(\left\{ f(\mathbf{x})^\top [f(\mathbf{x}^+) - f(\mathbf{x}_i^-)] \right\}_{i=1}^k \right) \right]. \quad (1)$$

Contrastive Learning Overview

We are interested in the excess risk $ER_{\text{un}}(\hat{f}_n) = L_{\text{un}}(\hat{f}_n) - \inf_{f \in \mathcal{F}} L_{\text{un}}(f)$ where \hat{f}_n is an empirical risk minimizer:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{j=1}^n \ell \left(\left\{ f(\mathbf{x}_j)^\top \left[f(\mathbf{x}_j^+) - f(\mathbf{x}_{ji}^-) \right] \right\}_{i=1}^k \right). \quad (2)$$

Previous works often assume that the tuples $\left\{ (\mathbf{x}_j, \mathbf{x}_j^+, \mathbf{x}_{j1:k}^-) \right\}_{j=1}^n$ are i.i.d. - A condition rarely met in practice.

Reality: Tuples are sub-sampled from a labeled dataset $\mathcal{S} = \left\{ (\mathbf{x}_j, \mathbf{y}_j) \right\}_{j=1}^N$.

Contrastive Learning Overview

In this work, we are interested in the excess risk bounds of $\hat{f}_{\text{sub}} = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}(f; \mathcal{T}_{\text{sub}})$ where:

$$\hat{\mathcal{L}}(f; \mathcal{T}_{\text{sub}}) = \frac{1}{n} \sum_{j=1}^n \ell \left(\left\{ f(\mathbf{z}_j)^\top \left[f(\mathbf{z}_j^+) - f(\mathbf{z}_{ji}^-) \right] \right\}_{i=1}^k \right). \quad (3)$$

Where $\mathcal{T}_{\text{sub}} = \left\{ (\mathbf{z}_j, \mathbf{z}_j^+, \mathbf{z}_{j1:k}^-) \right\}_{j=1}^M$ and each tuple is selected independently as:

- Select $c \in \mathcal{C}$ with probability $\hat{\rho}(c) = \frac{N_c^+}{N}$.
- Select $\mathbf{z}_j, \mathbf{z}_j^+$ from class c without replacement.
- Select $\mathbf{z}_{j1}^-, \dots, \mathbf{z}_{jk}^-$ outside of class c without replacement.

U-Statistics Formulation

For each $c \in \mathcal{C}$ let \mathcal{S}_c be the set of data points from class c and $\bar{\mathcal{S}}_c = \mathcal{S} \setminus \mathcal{S}_c$. Then, we define the set of $(k+2)$ -tuples \mathcal{T}_c as follows:

$$\mathcal{T}_c = \left\{ (\mathbf{z}, \mathbf{z}^+, \mathbf{z}_{1:k}^-) : \mathbf{z}, \mathbf{z}^+ \in \mathcal{S}_c, \mathbf{z}_1^-, \dots, \mathbf{z}_k^- \in \bar{\mathcal{S}}_c \right\}. \quad (4)$$

Then, we use the following U-Statistics to estimate $L_{\text{un}}(f)$:

$$\mathcal{U}_N(f) = \sum_{c \in \mathcal{C}} \frac{N_c^+}{N} \mathcal{U}_N(f|c), \quad (5)$$

$$\mathcal{U}_N(f|c) = \frac{1}{\binom{N_c^+}{2} \times \binom{N - N_c^+}{k}} \sum_{(\mathbf{z}, \mathbf{z}^+, \mathbf{z}_{1:k}^-) \in \mathcal{T}_c} \ell \left(\left\{ f(\mathbf{z})^\top [f(\mathbf{z}^+) - f(\mathbf{z}_i^-)] \right\}_{i=1}^k \right). \quad (6)$$

Contributions

Our contributions are as follows:

- We derive excess risk bounds for $\hat{f}_{\mathcal{U}} = \arg \min_{f \in \mathcal{F}} \mathcal{U}_N(f)$ and $\hat{f}_{\text{sub}} = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}(f; \mathcal{T}_{\text{sub}})$.
- We apply our bounds to linear maps and neural networks:

$$\mathcal{F}_{\text{lin}} = \left\{ \mathbf{x} \mapsto \mathbf{A}\mathbf{x} : \mathbf{A} \in \mathbb{R}^{d \times m}, \|\mathbf{A}^\top\|_{2,1} \leq a, \|\mathbf{A}\|_\sigma \leq s \right\}, \quad (7)$$

$$\mathcal{F}_{\text{nn}} = \mathcal{F}_L \circ \mathcal{F}_{L-1} \circ \cdots \circ \mathcal{F}_1, \quad (8)$$

$$\text{where } \mathcal{F}_l = \left\{ \mathbf{x} \mapsto \varphi_l(\mathbf{A}^{(l)}\mathbf{x}) : \mathbf{A}^{(l)} \in \mathbb{R}^{d_{l-1} \times d_l}, \|\mathbf{A}^{(l)}\|_\sigma \leq s_l \right\}.$$

Assumption: We assume that the contrastive loss $\ell : \mathbb{R}^k \rightarrow [0, \mathcal{M}]$ is ℓ^∞ -Lipschitz with constant $\eta > 0$.

Main Results

- **Theorem 5.1:** Let $\hat{f}_{\mathcal{U}} = \arg \min_{f \in \mathcal{F}} \mathcal{U}_N(f)$. For any $\delta \in (0, 1)$, with probability of at least $1 - \delta$:

$$\text{ER}_{\text{un}}(\hat{f}_{\mathcal{U}}) \leq \mathcal{O} \left(\sum_{c \in \mathcal{C}} \rho(c) \frac{K_{\mathcal{F},c}}{\sqrt{\tilde{N}}} + \mathcal{M} \sqrt{\frac{\ln |\mathcal{C}| / \delta}{\tilde{N}}} \right)$$

where $\tilde{N} = N \min \left(\frac{\rho_{\min}}{2}, \frac{1 - \rho_{\max}}{k} \right)$ and each term $K_{\mathcal{F},c}$ involves the covering number of \mathcal{F} .

- **Applications:** Bounds for linear maps and neural networks.

| Minimizer | Function Class | Generalization Bound |
|-------------------------|----------------|---|
| $\hat{f}_{\mathcal{U}}$ | Linear Maps | $\tilde{\mathcal{O}} \left(\frac{\eta s a b^2}{\sqrt{\tilde{N}}} + \mathcal{M} \sqrt{\frac{\ln \mathcal{C} / \delta}{\tilde{N}}} \right)$ |
| $\hat{f}_{\mathcal{U}}$ | Neural Nets | $\tilde{\mathcal{O}} \left(\frac{\mathcal{M} \mathcal{W}^{\frac{1}{2}}}{\sqrt{\tilde{N}}} + \mathcal{M} \sqrt{\frac{\ln \mathcal{C} / \delta}{\tilde{N}}} \right)$ |

Table: Generalization bounds for the U-Statistics minimizer.

Main Results

- **Theorem 5.2:** Let $\hat{f}_{\text{sub}} = \arg \min_{f \in \mathcal{F}} \hat{\mathcal{L}}(f; \mathcal{T}_{\text{sub}})$. For any $\delta \in (0, 1)$, with probability of at least $1 - \delta$:

$$\text{ER}_{\text{un}}(\hat{f}_{\text{sub}}) \leq \mathcal{O} \left(\hat{\mathfrak{R}}_{\mathcal{T}_{\text{sub}}}(\ell \circ \mathcal{F}) + \sum_{c \in \mathcal{C}} \rho(c) \frac{K_{\mathcal{F},c}}{\sqrt{\tilde{N}}} + \mathcal{M} \left[\sqrt{\frac{\ln |\mathcal{C}|/\delta}{\tilde{N}}} + \sqrt{\frac{\ln 1/\delta}{M}} \right] \right)$$

where $\tilde{N} = N \min \left(\frac{\rho_{\min}}{2}, \frac{1-\rho_{\max}}{k} \right)$ and each term $K_{\mathcal{F},c}$ involves the covering number of \mathcal{F} .

- **Applications:** Bounds for linear maps and neural networks.

| Minimizer | Function Class | Generalization Bound |
|------------------------|----------------|--|
| \hat{f}_{sub} | Linear Maps | $\tilde{\mathcal{O}} \left(\eta \text{sab}^2 \left[\frac{1}{\sqrt{\tilde{N}}} + \frac{1}{\sqrt{M}} \right] + \mathcal{M} \left[\sqrt{\frac{\ln \mathcal{C} /\delta}{\tilde{N}}} + \sqrt{\frac{\ln 1/\delta}{M}} \right] \right)$ |
| \hat{f}_{sub} | Neural Nets | $\tilde{\mathcal{O}} \left(\mathcal{M} \mathcal{W}^{\frac{1}{2}} \left[\frac{1}{\sqrt{\tilde{N}}} + \frac{1}{\sqrt{M}} \right] + \mathcal{M} \left[\sqrt{\frac{\ln \mathcal{C} /\delta}{\tilde{N}}} + \sqrt{\frac{\ln 1/\delta}{M}} \right] \right)$ |

Table: Generalization bounds for the sub-sampled empirical risk minimizer.