



## Contrastive learning for unsupervised classification

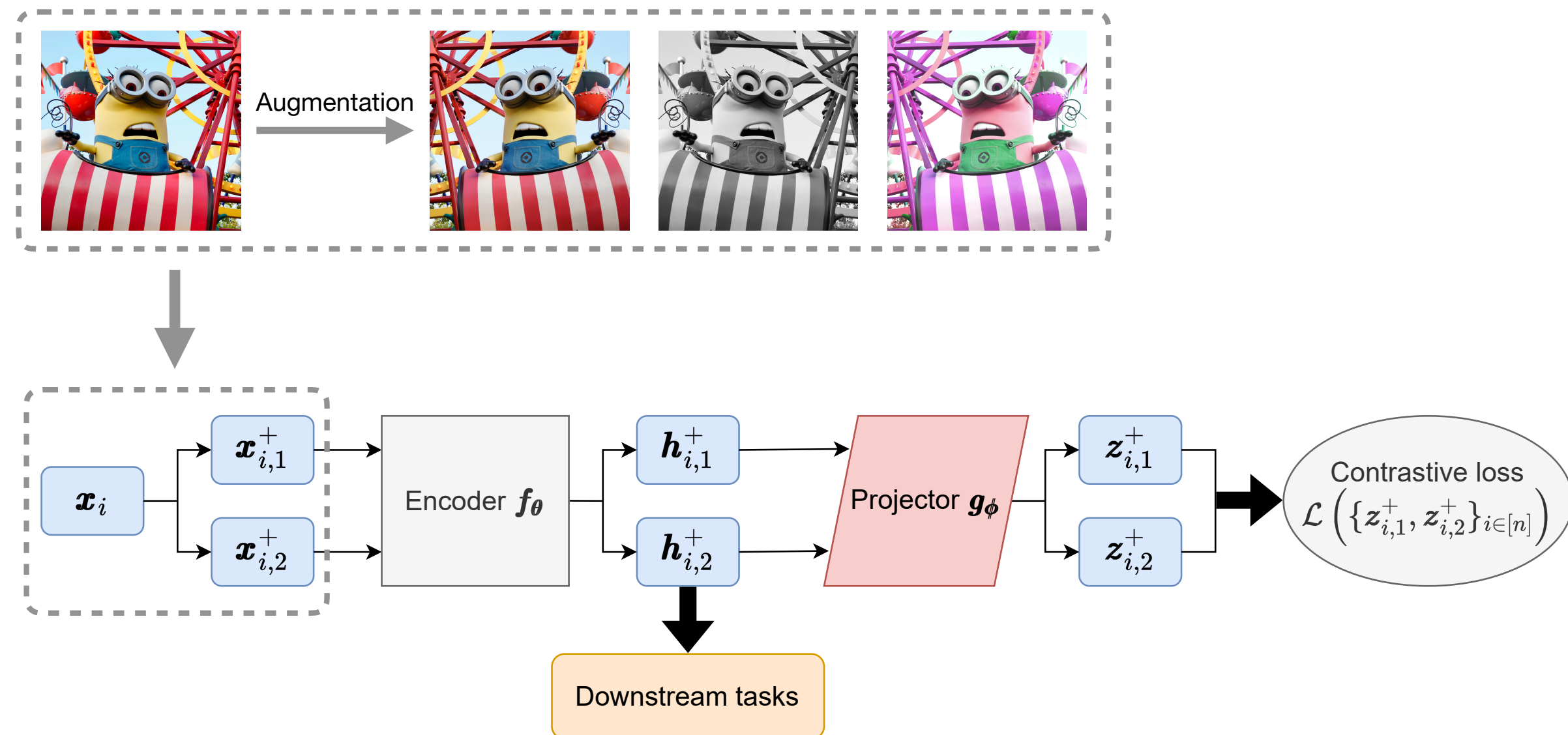


Figure 1. Encoder-projector framework.

- **Positive pairs:**  $(\mathbf{x}_{i,k}^+, \mathbf{x}_{i,l}^+)$ ; **Negative pairs:**  $(\mathbf{x}_{i,k}^+, \mathbf{x}_{j,l}^+)$ ,  $i \neq j$
- **Goal:** learn representations by
  1. encouraging proximity between positive pairs
  2. forcing negative pairs to be far
- Contrastive loss (cross-entropy with pseudo labels)

$$\min_{\theta, \varphi} \sum_i \sum_{j \sim i} -\log \frac{\exp(\frac{1}{\tau} \text{sim}(\mathbf{z}_i, \mathbf{z}_j))}{\sum_{k \neq i} \exp(\frac{1}{\tau} \text{sim}(\mathbf{z}_i, \mathbf{z}_k))}, \quad \text{sim}(\mathbf{z}, \mathbf{z}') = \frac{\langle \mathbf{z}, \mathbf{z}' \rangle}{\|\mathbf{z}\| \|\mathbf{z}'\|}$$

## Motivating questions

1. Effect of contrastive learning on representation and role of hyperparameters?
2. Causes of dimensional collapse in both features and embeddings?
3. Role of the projector? (removed after training in practice)

## Expansion and shrinkage of the signal

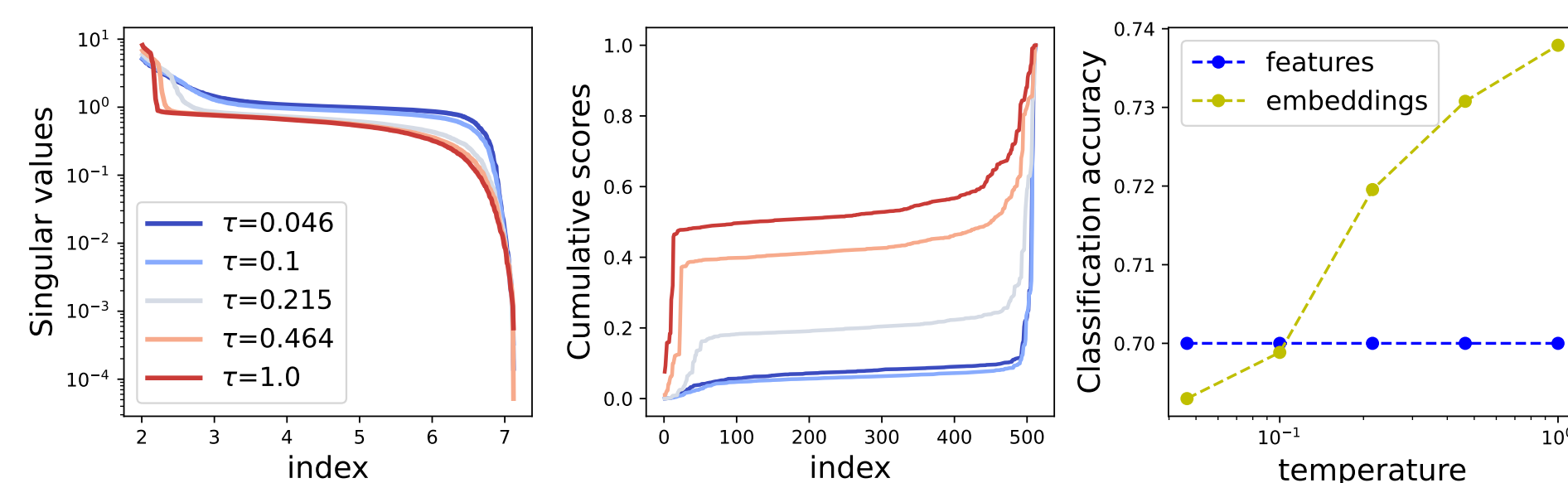


Figure 2. Results with the pretrained encoder and a one-layer linear projector.

$$\text{score}_i = \sum_{j \leq i} \frac{\langle \mathbf{v}_j, \boldsymbol{\mu}_{c_1, c_2} \rangle^2}{\|\boldsymbol{\mu}_{c_1, c_2}\|^2}$$

## Feature-level GMM modeling

- 2-GMM features:  $\mathbf{h}_i \stackrel{\text{i.i.d.}}{\sim} \frac{1}{2} \mathcal{N}(-\boldsymbol{\mu}, \mathbf{I}_p) + \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_p)$
- Augmentations:  $\mathbf{h}_{i,1}^+, \mathbf{h}_{i,2}^+ | \mathbf{h}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{h}_i, \sigma_{\text{aug}}^2 \mathbf{I}_p)$ ,  $\mathbf{h}_i^- \stackrel{d}{=} \mathbf{h}_{j,1}^+$ ,  $i \neq j$
- Linear projector:  $\mathbf{z}_i = \mathbf{W} \mathbf{h}_i$
- Population loss:

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2\tau} \cdot \frac{\mathbb{E}[\|\mathbf{W} \mathbf{h}_1^+ - \mathbf{W} \mathbf{h}_2^+\|^2]}{(\mathbb{E}[\|\mathbf{W} \mathbf{h}_1^+\|^2] \cdot \mathbb{E}[\|\mathbf{W} \mathbf{h}_2^+\|^2])^{1/2}} + \log \left( \mathbb{E} \exp \left( -\frac{1}{2\tau} \cdot \frac{\|\mathbf{W} \mathbf{h}_1^+ - \mathbf{W} \mathbf{h}_2^+\|^2}{(\mathbb{E}[\|\mathbf{W} \mathbf{h}_1^+\|^2] \cdot \mathbb{E}[\|\mathbf{W} \mathbf{h}_2^+\|^2])^{1/2}} \right) \right)$$

## Expansion-shrinkage phase transition in GMM features

Denote  $\tau^* = 2\|\boldsymbol{\mu}\|^2 \{(1 + \sigma_{\text{aug}}^2 + \|\boldsymbol{\mu}\|^2) \log(1 + 2\sigma_{\text{aug}}^2)\}^{-1}$ . A three-parameter configuration  $(\sigma_{\text{aug}}^2, \tau, \|\boldsymbol{\mu}\|^2)$  is said to be in the

- **expansion regime** if  $\tau \geq \tau^*$  and **shrinkage regime** if  $\tau < \tau^*$ .

## Theorem 1

Consider minimizer  $\mathbf{W}^*$  of certain first-order approximation  $\tilde{\mathcal{L}}(\mathbf{W})$ .

- When  $\tau \geq \tau^*$  (expansion regime),  $\mathbf{W}^* = \sum_j \sigma_j^* \mathbf{u}_j^* \mathbf{v}_j^{*\top}$  satisfies
 
$$\sigma_2^* = \dots = \sigma_p^* = 0, \quad \langle \mathbf{v}_1^*, \boldsymbol{\mu} \rangle^2 = \|\boldsymbol{\mu}\|^2 \quad \text{i.e., perfect alignment}$$
- When  $\tau < \tau^*$  (shrinkage regime),
 
$$\text{if } \sigma_{\text{aug}}^2 \rightarrow 0, \quad \text{then } \max_j |\sigma_j \langle \mathbf{v}_j^*, \boldsymbol{\mu} \rangle| \rightarrow 0 \quad \text{i.e., compress if correlated}$$

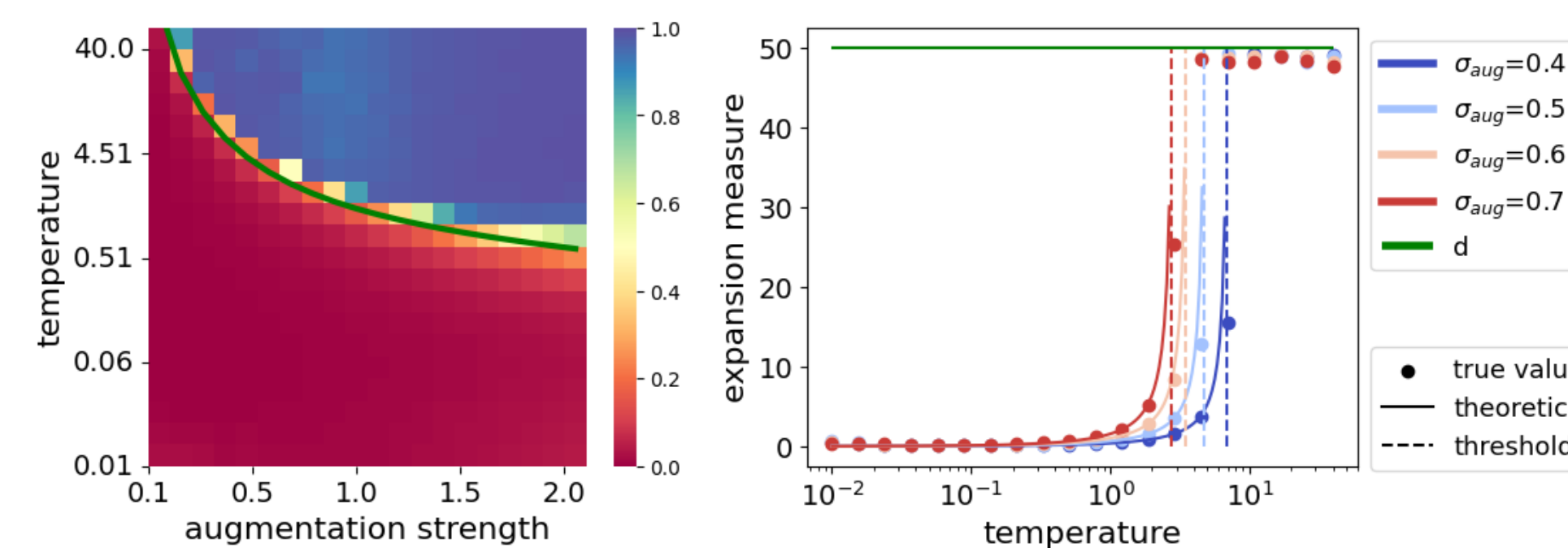


Figure 3. Expansion measure  $\tilde{r}(\mathbf{W}) = \|\mathbf{W} \boldsymbol{\mu}\|^2 / (\|\mathbf{W}\|_F^2 \|\boldsymbol{\mu}\|^2)$ .

## Empirical evidence for feature-level modeling

1. Linear separable features after a few epochs
2. Contrastive loss decomposition at each epoch  $t$ ,

$$\mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi}^{(t)}) = \min_{\boldsymbol{\varphi}} \mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi}) + \mathcal{L}^\perp(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi}^{(t)}),$$

which satisfies  $\mathcal{L}^\perp(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi}^{(t)}) \ll \min_{\boldsymbol{\varphi}} \mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi})$  and

$$\|\tilde{\boldsymbol{\varphi}}^{(t)} - \boldsymbol{\varphi}^{(t)}\| \ll \|\boldsymbol{\varphi}^{(t)}\|, \quad \tilde{\boldsymbol{\varphi}}^{(t)} = \text{argmin}_{\boldsymbol{\varphi}} \mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\varphi})$$

## Effect of projectors on generalization accuracy

- **Motivating question:** how does expansion/shrinkage affect generalization in downstream tasks?
- Consider the linear projection from the following class:

$$\mathcal{W} = \{\mathbf{W}_\eta = \mathbf{I}_p + \eta \cdot \rho^{-1} \boldsymbol{\mu} \boldsymbol{\mu}^\top : \eta > -1\}, \quad \text{where } \rho = \|\boldsymbol{\mu}\|^2$$

## Invariance of generalization error in low-dimensional regime

- $\ell_2$ -regularized logistic regression for non-separable data
 
$$\ell_n(\gamma, \boldsymbol{\beta}; \lambda_n) = \mathbb{E}_n \left\{ \log \left[ 1 + e^{-y(\gamma + \mathbf{z}^\top \boldsymbol{\beta})} \right] \right\} + \lambda_n \|\boldsymbol{\beta}\|^2$$
- Classification error  $\text{Err}(\gamma, \boldsymbol{\beta}; \eta, \lambda_n) = \mathbb{P}(\gamma + y'(\mathbf{z}', \boldsymbol{\beta}) < 0)$

## Proposition

Let  $(\hat{\gamma}, \hat{\boldsymbol{\beta}})$  be the minimizer of  $\ell_n(\gamma, \boldsymbol{\beta}; \lambda_n)$ .

1. If  $\lambda_n = a \cdot b_n > 0$  with constant  $a > 0$  and  $0 < b_n \ll \sqrt{n}$ , then the test error
 
$$\text{Err}(\hat{\gamma}, \hat{\boldsymbol{\beta}}; \eta, \lambda_n) = \Phi(-\|\boldsymbol{\mu}\|) + O_{\mathbb{P}}(b_n n^{-1/2}).$$
2. If  $\lambda_n = a\sqrt{n}$  with constant  $a > 0$ , then  $\text{Err}(\hat{\gamma}, \hat{\boldsymbol{\beta}}; \eta, \lambda_n)$  is decreasing in  $\eta$ .

## Decreasing generalization error in high-dimensional regime

- Implicit bias in *overparametrized* models: GD for logistic regression converges to max-margin classifier for separable data

$$\max_{\boldsymbol{\beta}} \min_{i \leq n} y_i \langle \mathbf{z}_i, \boldsymbol{\beta} \rangle$$

subject to  $\|\boldsymbol{\beta}\| \leq 1$

- Classification error  $\text{Err}(\hat{\boldsymbol{\beta}}; \eta) = \mathbb{P}(y'(\mathbf{z}', \hat{\boldsymbol{\beta}}) < 0)$
- A linear layer  $\mathbf{z}_i = \mathbf{W} \mathbf{h}_i$  can be interpreted as *reparametrization*

## Theorem 2

Suppose  $n/p \rightarrow \delta > 0$ . There exists threshold  $\delta^*(\rho) > 0$  such that

- (separability) if  $\delta < \delta^*$ , there exists a unique solution  $\hat{\boldsymbol{\beta}}$  with the margin

$$\hat{\kappa} = \min_{i \leq n} y_i \langle \mathbf{z}_i, \hat{\boldsymbol{\beta}} \rangle \xrightarrow{p} \kappa^*(\|\boldsymbol{\mu}\|, \eta) > 0$$

and conversely data are not separable w.h.p. if  $\delta > \delta^*$ .

- (monotone error) if  $\delta < \delta^*$ , the asymptotic error  $\text{Err}^*(\eta)$ , namely

$$\text{Err}(\hat{\boldsymbol{\beta}}; \eta) \xrightarrow{p} \text{Err}^*(\eta),$$

is decreasing in  $\eta$ .

## References

- [1] Yu Gui, Cong Ma, and Yiqiao Zhong. Demystifying projection heads in contrastive learning: an expansion and shrinkage perspective. *In preparation*, 2023.