

# A statistical theory of overfitting for imbalanced classification

Jingyang Lyu\*   Kangjie Zhou†   Yiqiao Zhong\*

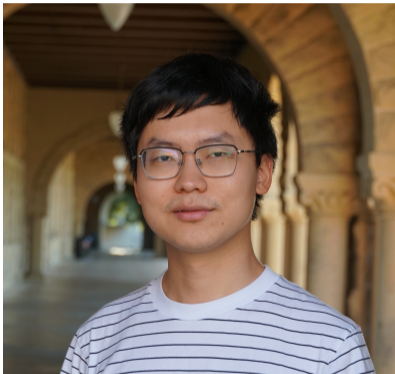
\*Department of Statistics, University of Wisconsin–Madison

†Department of Statistics, Columbia University



Department of Statistics  
SCHOOL OF COMPUTER, DATA & INFORMATION SCIENCES  
UNIVERSITY OF WISCONSIN-MADISON

# Collaborators



Kangjie Zhou, postdoc at Columbia U



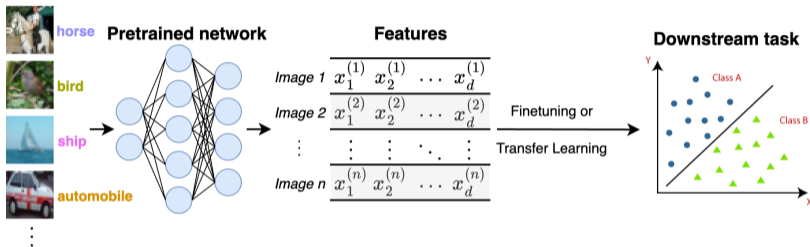
Yiqiao Zhong, UW–Madison

Paper: <https://arxiv.org/abs/2502.11323>

# Contents

- ▶ **Introduction**
- ▶ Settings
- ▶ Characterizing overfitting via empirical logit distribution
- ▶ Rebalancing margin is crucial
- ▶ Consequences for confidence estimation and calibration
- ▶ Generalization and future work

# Challenge 1: High dimensionality



High dimensional features are everywhere:

- Finetuning a classification layer in deep learning
- Linear probing, interpretability of LLMs
- Single-cell omics

# Challenge 1: High dimensionality

	Low dimensions	High dimensions
Parameter estimation	$\left\langle \frac{\hat{\beta}}{\ \hat{\beta}\ }, \frac{\beta}{\ \beta\ } \right\rangle \approx 1$	$\left\langle \frac{\hat{\beta}}{\ \hat{\beta}\ }, \frac{\beta}{\ \beta\ } \right\rangle < 1$
Generalization	Train error $\approx$ Test error	Train error $<$ Test error

**Table:** Qualitative comparison for linear classification,  $\beta$  is the slope parameter vector.

The advances of high-dimensional statistics in the past 15 years.

- El Karoui et al. (2013), Donoho and Montanari (2016), Sur and Candés (2019)
- Double descent and benign overfitting: Belkin et al. (2019), Bartlett et al. (2020)
- Many more ...

Q: New angles for the (overfitting) effects of dimensionality?

## Challenge 2: Data imbalance

Real-world datasets are generally **imbalanced**.

- **Sentiment analysis.**

Dataset	Tweets	#Negative	#Positive
<i>Stanford Twitter Test Set (STS-Test)</i> [9]	359	177	182
<i>Sanders Dataset (Sanders)</i> [17]	1224	654	570
<i>Obama McCain Debate (OMD)</i> [7]	1906	1196	710
<i>Health Care Reform (HCR)</i> [22]	1922	1381	541
<i>Stanford Gold Standard (STS-Gold)</i> [17]	2034	632	1402
<i>Sentiment Strength Twitter Dataset (SSTD)</i> [23]	2289	1037	1252
<i>The Dialogue Earth Weather Dataset (WAB)</i> [3]	5495	2580	2915
<i>The Dialogue Earth Gas Prices Dataset (GASP)</i> [3]	6285	5235	1050
<i>Semeval Dataset (Semeval)</i> [14]	7535	2186	5349

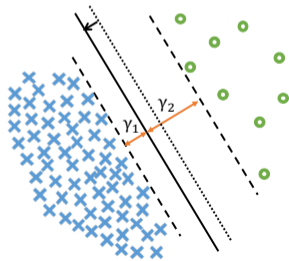
**Figure:** Twitter datasets used for sentiment analysis [Saif et al. 2015]

- **Industrial fault detection** (failures  $\ll$  normal operations)
- **Healthcare and medical diagnosis** (rare disease/genetic markers, privacy issue)

## Challenge 2: Data imbalance

- **Minority classes** have worse training and testing errors.
- The classical asymptotic theory or finite-sample analysis is inaccurate in high dimensions.
- The practice is heuristic-driven and ad hoc.
  - Re-sampling: oversampling the minority or under-sampling the majority
  - Re-weighting: assigning higher weights for minority classes
  - Synthetic data: SMOTE (2002), Mixup (2018)
  - **Margin adjustment**: popular in deep learning.

Q: How to quantify the impact of factors  
(imbalance ratio, SNR, dimension) on accuracy?



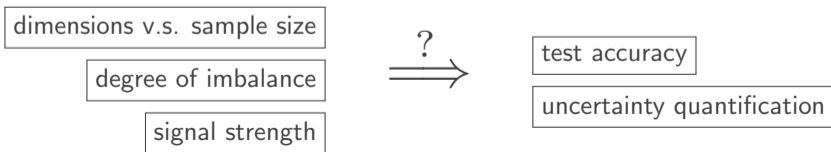
Cao et al. 2019

# Goals of this talk

**Goal 1.** Provide a new angle of **characterizing overfitting** for imbalanced classification.

	Low dimensions	High dimensions
Parameter estimation	$\left\langle \frac{\hat{\beta}}{\ \hat{\beta}\ }, \frac{\beta}{\ \beta\ } \right\rangle \approx 1$	$\left\langle \frac{\hat{\beta}}{\ \hat{\beta}\ }, \frac{\beta}{\ \beta\ } \right\rangle < 1$
Generalization	Train error $\approx$ Test error	Train error $<$ Test error
Distribution of logits	1D projection of $P_x$	Skewed/distorted 1D projection of $P_x$

**Goal 2.** Quantify the **adverse effects** of overfitting, esp. for the minority class.



# Contents

- ▶ Introduction
- ▶ **Settings**
- ▶ Characterizing overfitting via empirical logit distribution
- ▶ Rebalancing margin is crucial
- ▶ Consequences for confidence estimation and calibration
- ▶ Generalization and future work

# Binary classification

- Training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \stackrel{\text{i.i.d.}}{\sim} P_{\mathbf{x}, y}$ .
  - $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{+1, -1\}$
- **Imbalance ratio:** denote  $\pi = \mathbb{P}(y_i = +1)$ .
  - The classification is imbalanced if  $\pi < 1/2$ .

$$y_i = \begin{cases} +1 & \text{with prob } \pi \quad (\text{minority}) \\ -1 & \text{with prob } 1 - \pi \quad (\text{majority}) \end{cases}$$

- Build a classifier based on  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ .

For a point  $\mathbf{x}$ , the predicted label is

$$\hat{y}(\mathbf{x}) = \begin{cases} +1 & \text{if } f(\mathbf{x}) > 0 \\ -1 & \text{if } f(\mathbf{x}) \leq 0 \end{cases}$$

# Two linear classifiers

- We focus on two linear classifiers.

$$\text{(logistic regression)} \quad \underset{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i(\langle \mathbf{x}_i, \beta \rangle + \beta_0)),$$

$$\begin{aligned} \text{(SVM)} \quad & \underset{\beta \in \mathbb{R}^d, \beta_0, \kappa \in \mathbb{R}}{\text{maximize}} \quad \kappa, \\ & \text{subject to} \quad y_i(\langle \mathbf{x}_i, \beta \rangle + \beta_0) \geq \kappa, \quad \forall 1 \leq i \leq n, \\ & \quad \|\beta\|_2 \leq 1. \end{aligned}$$

- **Connection** (inductive bias): when training data is linear separable,

$$\text{SVM} = \text{Max-margin classifier} = \text{Ridgeless logistic regression}$$

# Contents

- ▶ Introduction
- ▶ Settings
- ▶ **Characterizing overfitting via empirical logit distribution**
- ▶ Rebalancing margin is crucial
- ▶ Consequences for confidence estimation and calibration
- ▶ Generalization and future work

# Logits and empirical logit distribution

For any classifier  $\hat{y}(x) = 2\mathbb{1}\{\hat{f}(x) > 0\} - 1$  (e.g., SVM, neural network, language model)

- **Logit** for point  $x$ :  $\hat{f}(x)$
- **Margin**:  $\hat{\kappa}_n = \min_{1 \leq i \leq n} y_i \hat{f}(x_i)$ 
  - When  $\hat{\kappa}_n > 0$ , the training set is **linearly separable**.

## Definition (Empirical logit distribution, or ELD)

For any binary classifier  $\hat{y}(x)$  built on  $\hat{f}(x)$ , the *empirical logit distribution* is defined as

$$\hat{\nu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{(y_i, \hat{f}(x_i))} \quad (2)$$

where  $\delta_a$  denotes the delta measure supported at point  $a$ .

# Empirical logit distribution v.s. testing logit distribution

Let  $(\mathbf{x}_{\text{test}}, y_{\text{test}}) \sim P_{\mathbf{x}, y}$  be a new data point.

- **Overfitting** can be characterized by discrepancy between

$$\underbrace{\hat{\nu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{(y_i, \hat{f}(\mathbf{x}_i))}}_{\substack{\text{empirical logit distribution} \\ \text{("training" logit distribution)}}} \quad \text{and} \quad \underbrace{\hat{\nu}_n^{\text{test}} = \text{Law}(y_{\text{test}}, \hat{f}(\mathbf{x}_{\text{test}}))}_{\text{testing logit distribution}}$$

- Note: both  $\hat{\nu}_n, \hat{\nu}_n^{\text{test}}$  are random measures.
  - Since  $\hat{f}$  depends on training set  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .

# Empirical phenomenon: Simulation

Settings:

1. Generate a **(linearly) separable** training set from a Gaussian mixture model (GMM):

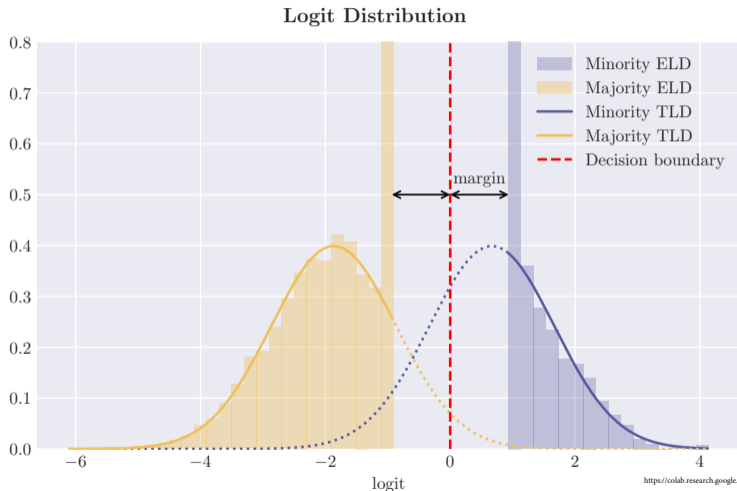
$$y_i = \begin{cases} +1, & \text{w.p. } \pi \quad (\text{minority}) \\ -1, & \text{w.p. } 1 - \pi \quad (\text{majority}) \end{cases}, \quad \mathbf{x}_i | y_i \sim \mathcal{N}(y_i \boldsymbol{\mu}, \mathbf{I}_d), \quad i = 1, 2, \dots, n.$$

2. Train a **max-margin classifier (SVM)**:  $\implies \hat{\boldsymbol{\beta}}, \hat{\beta}_0, \hat{\kappa}$

$$\begin{aligned} & \underset{\boldsymbol{\beta} \in \mathbb{R}^d, \beta_0 \in \mathbb{R}, \kappa \in \mathbb{R}}{\text{maximize}} && \kappa, \\ & \text{subject to} && y_i(\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \beta_0) \geq \kappa, \quad \forall 1 \leq i \leq n, \\ & && \|\boldsymbol{\beta}\|_2 \leq 1. \end{aligned}$$

3. Compare empirical / testing logit distribution for  $\hat{f}(\mathbf{x}) = \langle \mathbf{x}, \hat{\boldsymbol{\beta}} \rangle + \hat{\beta}_0$ .

# Empirical phenomenon: Simulation



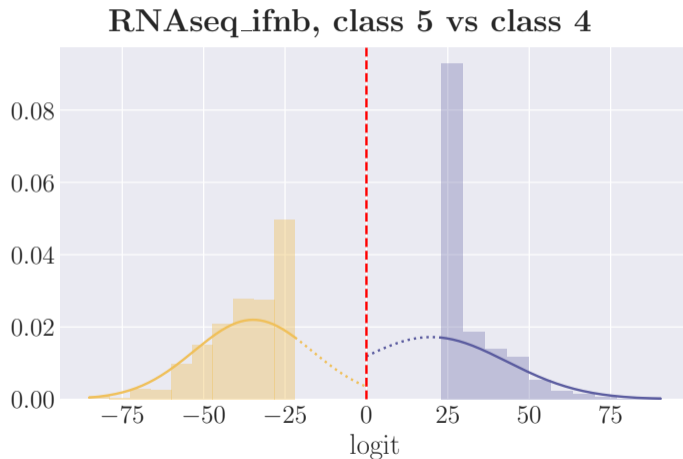
Code



Figure: Empirical (training) and testing logit distribution for binary Gaussian mixture model

# Empirical phenomenon: tubular data

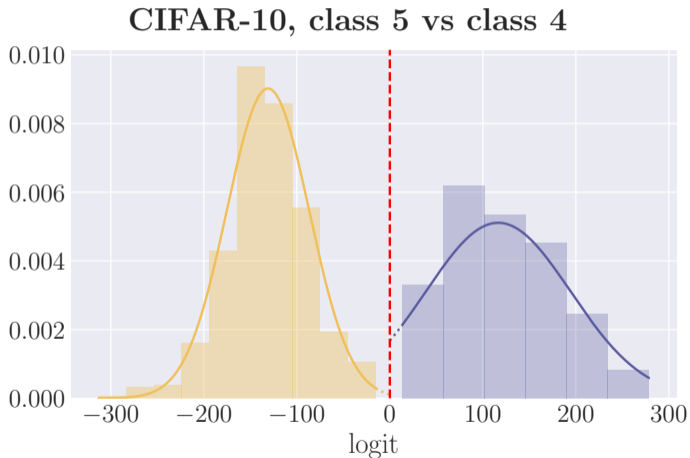
RNA-seq ifnb dataset with logistic regression ( $\pi = 0.2$ )



**Figure:** Empirical (training) and testing logit distribution for single-cell dataset

# Empirical phenomenon: image data

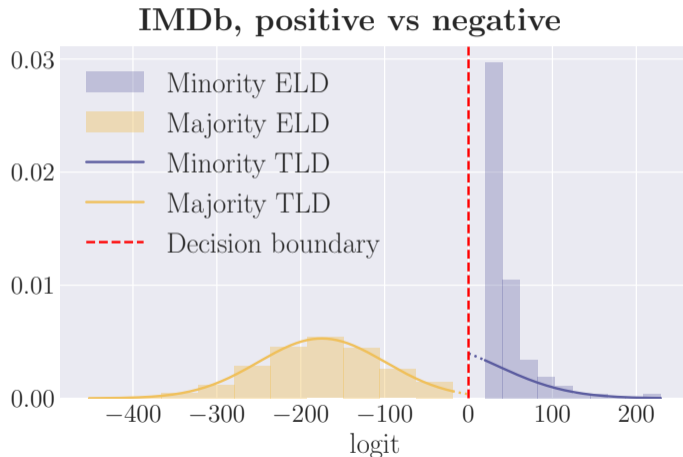
ResNet-18 trained on CIFAR-10 ( $\pi = 0.1$ )



**Figure:** Empirical (training) and testing logit distribution for CIFAR-10 dataset

# Empirical phenomenon: text data

BERT(110M) trained on IMDb movie reviews ( $\pi = 0.02$ )



**Figure:** Empirical (training) and testing logit distribution for IMDb dataset

# Theoretical foundation

Consider GMM with asymptotic regime  $n/d \rightarrow \delta \in (0, \infty)$ .

- Recall  $\hat{\kappa} = \min_{1 \leq i \leq n} y_i (\langle \mathbf{x}_i, \hat{\boldsymbol{\beta}} \rangle + \hat{\beta}_0)$ . Denote  $\hat{\rho} = \left\langle \frac{\hat{\boldsymbol{\beta}}}{\|\hat{\boldsymbol{\beta}}\|}, \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|} \right\rangle$ . ( $\|\hat{\boldsymbol{\beta}}\| = 1$  when separable)
- We may expect  $(\hat{\rho}, \hat{\beta}_0, \hat{\kappa})$  converge to some limit  $(\rho^*, \beta_0^*, \kappa^*)$  as  $n, d \rightarrow \infty$ .

Let  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$  be a new testing point, then

$$\begin{aligned} y_{\text{test}} \left( \langle \mathbf{x}_{\text{test}}, \hat{\boldsymbol{\beta}} \rangle + \hat{\beta}_0 \right) &= y_{\text{test}} \left\langle y_{\text{test}} \boldsymbol{\mu} + \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \hat{\boldsymbol{\beta}} \right\rangle + y_{\text{test}} \hat{\beta}_0 \\ &= \hat{\rho} \|\boldsymbol{\mu}\| + \left\langle \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \hat{\boldsymbol{\beta}} \right\rangle + y_{\text{test}} \hat{\beta}_0 \\ &\approx \rho^* \|\boldsymbol{\mu}\| + G + Y \beta_0^*, \quad \text{where } (Y, G) \sim P_y \times \mathcal{N}(0, 1). \end{aligned}$$

# Theoretical foundation

For a testing point  $(\mathbf{x}_{\text{test}}, y_{\text{test}})$ ,

$$y_{\text{test}} \left( \langle \mathbf{x}_{\text{test}}, \hat{\boldsymbol{\beta}} \rangle + \hat{\beta}_0 \right) \approx \rho^* \|\boldsymbol{\mu}\| + G + Y \beta_0^*. \quad (\hat{\nu}_n^{\text{test}})$$

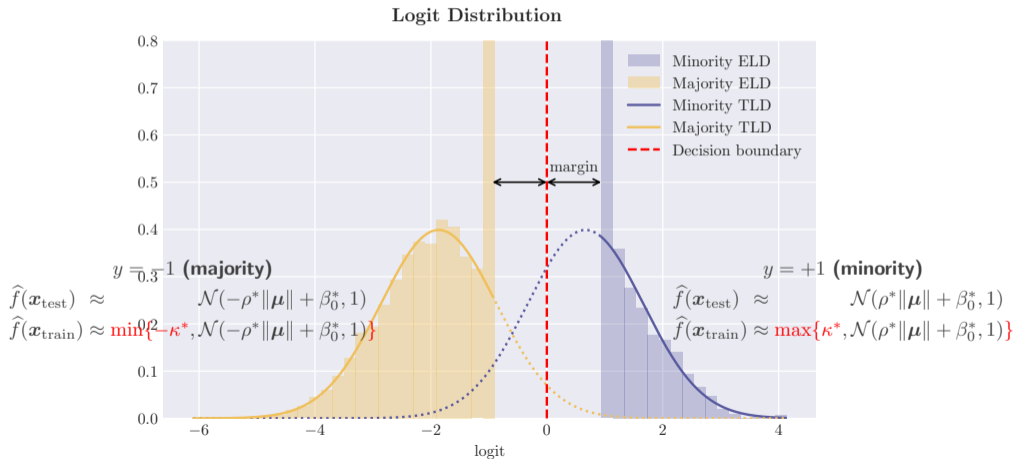
However, for a training point  $(\mathbf{x}_i, y_i)$ ,

- There is a **distortion effect** on the distribution due to dependence between  $(\mathbf{x}_i, y_i)$  and  $\hat{f}$ .

$$y_i \left( \langle \mathbf{x}_i, \hat{\boldsymbol{\beta}} \rangle + \hat{\beta}_0 \right) \approx \max \left\{ \kappa^*, \rho^* \|\boldsymbol{\mu}\| + G + Y \beta_0^* \right\}. \quad (\hat{\nu}_n)$$

# Key takeaway

## Overfitting = “Truncation”



# Theoretical foundation

## Theorem (Separable regime, simplified ver.)

Consider GMM with asymptotic regime  $n/d \rightarrow \delta \in (0, \infty)$ .

(a) **(Phase transition)** There is a critical threshold  $\delta_c = \delta_c(\|\boldsymbol{\mu}\|, \pi)$ , such that

$$\mathbb{P}\{\text{training set is linearly separable}\} \rightarrow 1, \quad \text{if } \delta < \delta_c.$$

(b) **(Parameter convergence)** If  $\delta < \delta_c$ , then  $(\hat{\rho}, \hat{\beta}_0, \hat{\kappa}) \xrightarrow{P} (\rho^*, \beta_0^*, \kappa^*)$ , where  $(\rho^*, \beta_0^*, \kappa^*)$  is the unique solution of the following variational optimization problem:

$$\begin{aligned} & \underset{\rho \in [-1, 1], \beta_0 \in \mathbb{R}, \kappa > 0, \boldsymbol{\xi} \in \mathcal{L}^2}{\text{maximize}} && \kappa, \\ & \text{subject to} && \rho \|\boldsymbol{\mu}\| + G + Y\beta_0 + \sqrt{1 - \rho^2} \boldsymbol{\xi} \geq \kappa, \quad \mathbb{E}[\boldsymbol{\xi}^2] \leq 1/\delta. \end{aligned}$$

(c) **(ELD convergence)** If  $\delta < \delta_c$ , denote  $\nu^* = \max\{\kappa^*, \rho^* \|\boldsymbol{\mu}\| + G + Y\beta_0^*\}$ . Then  $W_2(\hat{\nu}_n, \nu^*) \xrightarrow{P} 0$ .

# Theoretical foundation: remarks

---


$$\begin{array}{ll}
 \underset{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}, \kappa \in \mathbb{R}}{\text{maximize}} & \kappa, \\
 \text{subject to} & \forall 1 \leq i \leq n \quad y_i(\langle \mathbf{x}_i, \beta \rangle + \beta_0) \geq \kappa, \quad \|\beta\|_2 \leq 1. \quad (\text{A})
 \end{array}$$


---

$$\begin{array}{ll}
 \underset{\rho \in [-1, 1], \beta_0 \in \mathbb{R}, \kappa > 0, \xi \in \mathcal{L}^2}{\text{maximize}} & \kappa, \\
 \text{subject to} & \rho \|\mu\| + G + Y\beta_0 + \sqrt{1 - \rho^2} \xi \geq \kappa, \quad \mathbb{E}[\xi^2] \leq 1/\delta. \quad (\text{B})
 \end{array}$$


---

- In (B), it can be shown that  $\sqrt{1 - \rho^2} \xi = (\kappa - \rho \|\mu\| - G - Y\beta_0)_+ \quad (t)_+ = \max\{0, t\}$ .  
 $\Rightarrow$  The random variable  $\xi$  represents the **overfitting effect** in high dimensions.
- In (B),  $\beta_0^* < 0$ . The mean of minority testing logits is *closer to margin* than majority.  
 $\Rightarrow$  **Overfitting hurts minority** class more than majority.

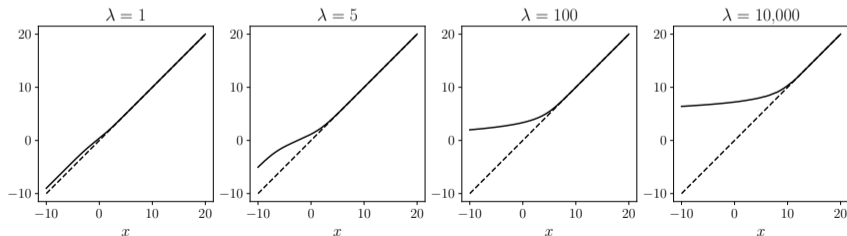
# Theoretical foundations: non-separable regime

**Logistic regression:** we obtained similar variational formulation in the limit.

$$\begin{aligned} & \underset{\rho \in [-1,1], R \geq 0, \beta_0 \in \mathbb{R}, \xi \in \mathcal{L}^2}{\text{minimize}} && \mathbb{E} \left[ \ell \left( \rho \|\boldsymbol{\mu}\|_2 R + RG + Y\beta_0 + R\sqrt{1 - \rho^2 \xi} \right) \right], \\ & \text{subject to} && \mathbb{E} [\xi^2] \leq 1/\delta. \end{aligned}$$

$\text{prox}_{\lambda \ell}(x) := \arg \min_{t \in \mathbb{R}} \{ \ell(t) + \frac{1}{2\lambda}(t - x)^2 \}$

**Proximal operator** instead of truncation characterizes overfitting effects.



**Figure:** Plots of proximal operator  $x \mapsto \text{prox}_{\lambda \ell}(x)$  where  $\lambda$  represents the strength of overfitting.

# Contents

- ▶ Introduction
- ▶ Settings
- ▶ Characterizing overfitting via empirical logit distribution
- ▶ **Rebalancing margin is crucial**
- ▶ Consequences for confidence estimation and calibration
- ▶ Generalization and future work

# Rebalancing margin

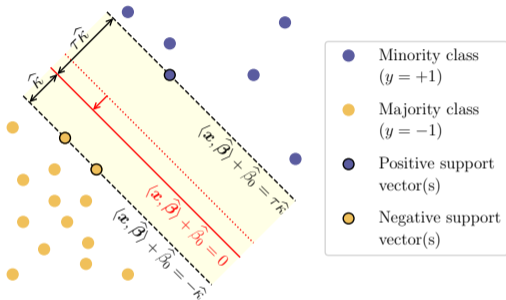
Rebalancing margin is crucial in separable regime.

Consider **margin-rebalanced SVM**:

$$\begin{aligned} & \underset{\beta \in \mathbb{R}^d, \beta_0 \in \mathbb{R}, \kappa \in \mathbb{R}}{\text{maximize}} && \kappa, \\ & \text{subject to} && y_i(\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \beta_0) \geq \tau \kappa, \quad \forall i : y_i = +1 \\ & && y_i(\langle \mathbf{x}_i, \boldsymbol{\beta} \rangle + \beta_0) \geq \kappa, \quad \forall i : y_i = -1 \\ & && \|\boldsymbol{\beta}\|_2 \leq 1. \end{aligned}$$

**Margin ratio:**  $\tau > 0$ .

- **Note:**  $\hat{\boldsymbol{\beta}}$  does not depend on  $\tau$ .
- Question: what is the optimal  $\tau$ ?



# Classification errors

	Exact testing error	Asymptotic testing error
	(minority error) $\text{Err}_+ = \mathbb{P}(\hat{y}(\mathbf{x}) \neq y \mid y = +1)$	$\rightarrow \text{Err}_+^* = \Phi(-\rho^* \ \boldsymbol{\mu}\  - \beta_0^*)$
	(majority error) $\text{Err}_- = \mathbb{P}(\hat{y}(\mathbf{x}) \neq y \mid y = -1)$	$\rightarrow \text{Err}_-^* = \Phi(-\rho^* \ \boldsymbol{\mu}\  + \beta_0^*)$
✗	(total error) $\text{Err} = \mathbb{P}(\hat{y}(\mathbf{x}) \neq y)$	$\rightarrow \text{Err}^* = \pi \text{Err}_+^* + (1 - \pi) \text{Err}_-^*$
✓	(balanced error) $\text{Err}_b = \frac{1}{2} \text{Err}_+ + \frac{1}{2} \text{Err}_-$	$\rightarrow \text{Err}_b^* = \frac{1}{2} \text{Err}_+^* + \frac{1}{2} \text{Err}_-^*$

# Setting 1: proportional regime

## Simulations

**Setup:** sample size  $n = 100$ , dimension  $d = 200$ . Run SVM, report errors over 100 runs.



Figure: Effects of margin rebalancing on test errors.

# Setting 1: proportional regime

## Simulations

**Setup:** sample size  $n = 100$ , dimension  $d = 200$ . Run SVM, report errors over 100 runs.

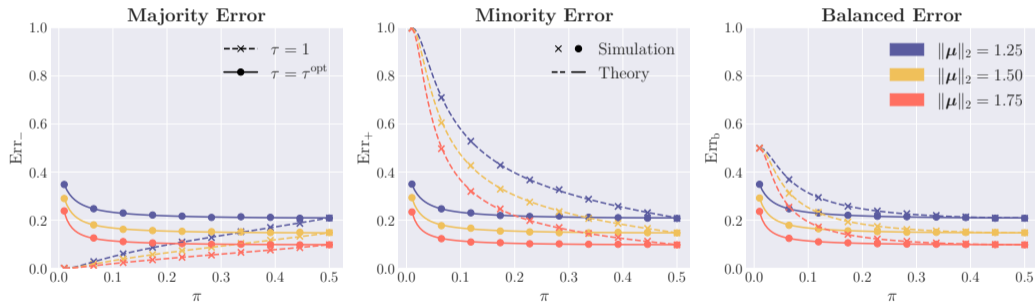


Figure: Impact of imbalance on test errors.

# Setting 1: proportional regime

Theoretical foundation

## Proposition (Proportional regime)

Define  $\tau^{\text{opt}}$  as the optimal margin ratio which minimizes the asymptotic balanced error

$$\tau^{\text{opt}} := \arg \min_{\tau \geq 1} \text{Err}_b^* = \arg \min_{\tau \geq 1} \{ \Phi(-\|\mu\|_2 \rho^* - \beta_0^*) + \Phi(-\|\mu\|_2 \rho^* + \beta_0^*) \}.$$

(a) When  $\tau = \tau^{\text{opt}}$ , we have  $\beta_0^* = 0$  and  $\text{Err}_+^* = \text{Err}_-^* = \text{Err}_b^*$ . In particular,

$$\tau^{\text{opt}} = \frac{g_1^{-1} \left( \frac{\rho^*}{2\pi \|\mu\|_2 \delta} \right) + \rho^* \|\mu\|_2}{g_1^{-1} \left( \frac{\rho^*}{2(1-\pi) \|\mu\|_2 \delta} \right) + \rho^* \|\mu\|_2}, \quad \text{where} \quad \begin{aligned} g_1(t) &= \mathbb{E}[(G + t)_+] \\ G &\sim \mathcal{N}(0, 1), (t)_+ = 0 \vee t \end{aligned}$$

(b) When  $\tau = \tau^{\text{opt}}$ , the testing error  $\text{Err}_b^*$  is a **decreasing** function of  $\|\mu\|_2$  (signal strength),  $\delta$  (aspect ratio) and  $\pi \in (0, 1/2)$  (imbalance ratio).

- When  $\pi$  is small, roughly speaking  $\tau^{\text{opt}} \asymp 1/\sqrt{\pi}$ .

## Setting 2: high imbalance

$$\pi \rightarrow 0, \|\boldsymbol{\mu}\| \rightarrow \infty, \delta = n/d \rightarrow \infty$$

- Motivation: in overparametrized model, the imbalance ratio ( $\pi$ ) is **vanishingly small** relative to dimension ( $d$ ) and sample size ( $n$ ).

Under Gaussian mixture model, consider ( $a, b, c > 0$ )

$$\pi \asymp d^{-a}, \quad \|\boldsymbol{\mu}\|^2 \asymp d^b, \quad n \asymp d^{c+1}.$$

- Such high imbalance dataset is **always separable** (with high probability).
- Feature distribution can be generalized to **sub-Gaussian**.

## Setting 2: high imbalance: phase transition

### Theorem (High imbalance regime, sub-Gaussian mixture model)

Suppose that  $a - c < 1$  (i.e.  $n\pi \rightarrow \infty$ ).

(a) **High signal** (no need for margin rebalancing):  $a - c < b$ . If  $1 \leq \tau_d \ll d^{b/2}$ , then

$$\text{Err}_+^* = o(1), \quad \text{Err}_-^* = o(1).$$

(b) **Moderate signal** (margin rebalancing is crucial):  $b < a - c < 2b$ . If we choose  $d^{a-b-c} \ll \tau_d \ll d^{(a-c)/2}$ , then

$$\text{Err}_+^* = o(1), \quad \text{Err}_-^* = o(1).$$

However, if we naively choose  $\tau_d \asymp 1$ , then

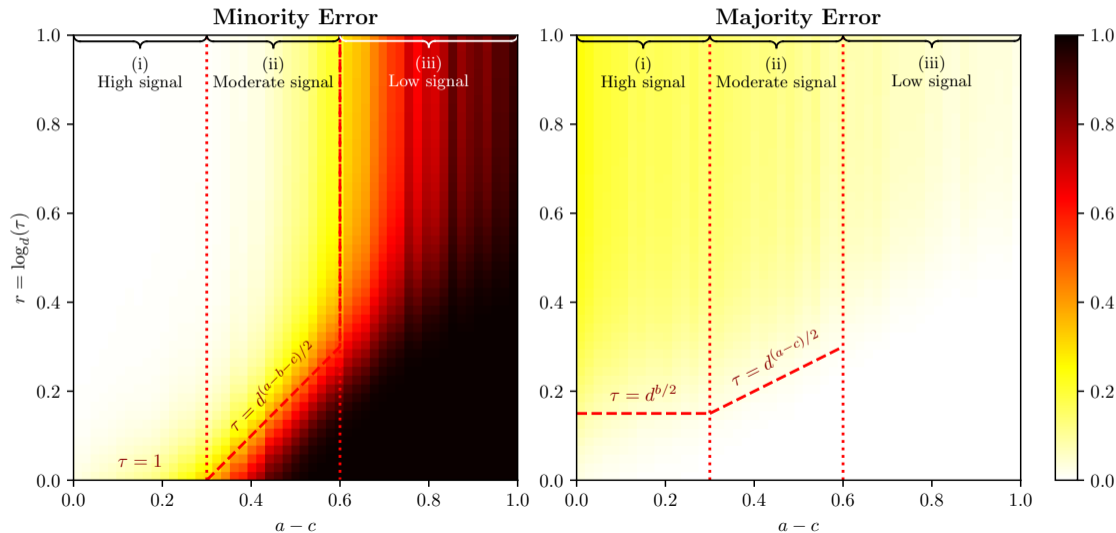
$$\text{Err}_+^* = 1 - o(1), \quad \text{Err}_-^* = o(1).$$

(c) **Low signal** (no better than random guess):  $a - c > 2b$ . For any  $\tau_d$ , we have

$$\text{Err}_b^* \geq \frac{1}{2} - o(1).$$

# Simulation: $\tau = d^r$

$\pi \asymp d^{-a}$ ,  $\|\mu\| \asymp d^{b/2}$ ,  $n \asymp d^{c+1}$  (fix  $b = 0.3$ ,  $c = 0.1$ ,  $d = 2000$ )



# Contents

- ▶ Introduction
- ▶ Settings
- ▶ Characterizing overfitting via empirical logit distribution
- ▶ Rebalancing margin is crucial
- ▶ **Consequences for confidence estimation and calibration**
- ▶ Generalization and future work

# Confidence and Calibration

**Confidence** (predicted probability)

- Multiclass classification: `softmax`
- Binary classification: sigmoid transformation  $p(\mathbf{x}) = 1/[1 + \exp(-f(\mathbf{x}))]$ .

Ideally, we expect  $p(\mathbf{x}) \approx \mathbb{P}(y = 1 \mid \mathbf{x})$ . But the RHS is often intractable in high dim.

## Definition (calibration)

A function  $p : \mathcal{X} \rightarrow [0, 1]$  is (*perfectly*) *calibrated* if

$$p(\mathbf{x}) = \mathbb{P}(y = 1 \mid p(\mathbf{x})) \quad \text{a.s.}$$

**Intuition:** Given 1,000 predictions, each with confidence of 0.2, we expect that about 200 should be classified as positive.

- Most informative example:  $p(\mathbf{x}) = \mathbb{P}(y = 1 \mid \mathbf{x})$ .
- Least informative example:  $p(\mathbf{x}) \equiv \mathbb{P}(y = 1) = \pi$ .

# Calibration and other uncertainty measurements

**Calibration error (CE).**

$$\text{CE}(p) = \mathbb{E} \left[ \left( \mathbb{P}(y = 1 \mid p(\mathbf{x})) - p(\mathbf{x}) \right)^2 \right]$$

- Calibration itself does not guarantee a useful predictor, e.g.,  $p(\mathbf{x}) = \pi$ .
- The variance in  $y$  explained by prediction  $p(\mathbf{x})$  shouldn't be too small (**Sharpness**).

**Mean squared error (MSE).**

$$\text{MSE}(p) = \mathbb{E} \left[ \left( \mathbb{1}\{y = 1\} - p(\mathbf{x}) \right)^2 \right]$$

**Confidence estimation error (ConfErr).**

$$\text{ConfErr}(p) := \mathbb{E} \left[ \left( \mathbb{P}(y = 1 \mid \mathbf{x}) - p(\mathbf{x}) \right)^2 \right] .$$

# Calibration: simulation

**Setup:** 2-GMM,  $n = 1,000$ ,  $d = 500$ ,  $\pi = 0.05$ ,  $\|\mu\| = 1$ , train SVM with  $\tau = \tau^{\text{opt}}$ .

**Reliability diagrams:** For each  $p$  ( $x$ -axis), calculate  $\mathbb{P}(y = 1 \mid \hat{p}(x) = p)$  ( $y$ -axis) on test set.

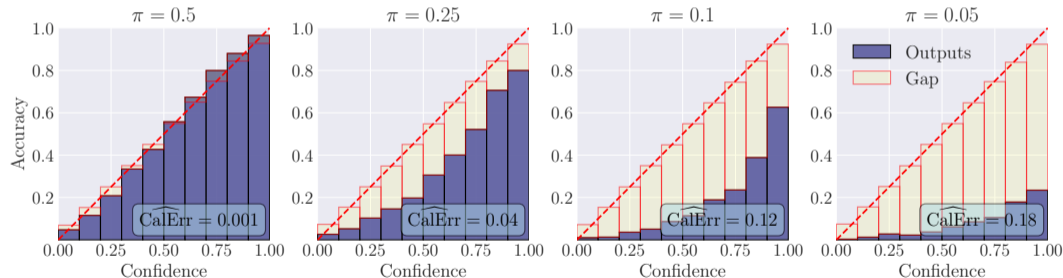


Figure: Imbalance worsens calibration.

# Confidence and calibration: Theoretical foundations

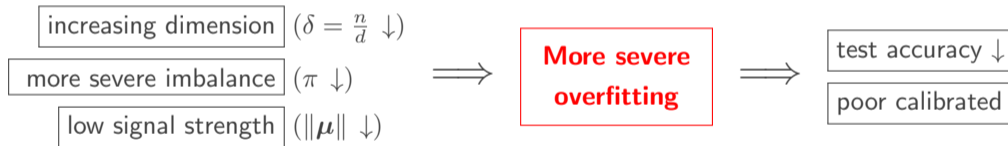
Under proportional regime  $n, d \rightarrow \infty$ ,  $n/d \rightarrow \delta$ , we show:

	$\text{Err}_+^*, \text{Err}_-^*, \text{Err}_b^*$	$\text{CE}^*$	$\text{MSE}^*$	$\text{ConfErr}^*$
imbalance ratio $\pi \uparrow$	$\downarrow$		$\downarrow$	$\downarrow$
signal strength $\ \boldsymbol{\mu}\ _2 \uparrow$	$\downarrow$	$\downarrow$	$\downarrow$	
aspect ratio $n/d \rightarrow \delta \uparrow$	$\downarrow$	$\downarrow$	$\downarrow$	$\downarrow$

**Table:** Monotonicity of test errors and confidence/calibration metrics

Qualitatively, the effects of imbalance is similar to  
signal strength and effective sample size.

# Takeaway message



# Contents

- ▶ Introduction
- ▶ Settings
- ▶ Characterizing overfitting via empirical logit distribution
- ▶ Rebalancing margin is crucial
- ▶ Consequences for confidence estimation and calibration
- ▶ **Generalization and future work**

# Generalization

- **Non-isotropic covariance.**
  - We obtained a variational form based on formal calculation.
  - Dependence on the covariance spike and direction is complicated.
- **Multiclass classification.**
  - Truncation for 2-dim Gaussian can be observed for empirical logit distribution.

# Multiclass classification: CIFAR-10

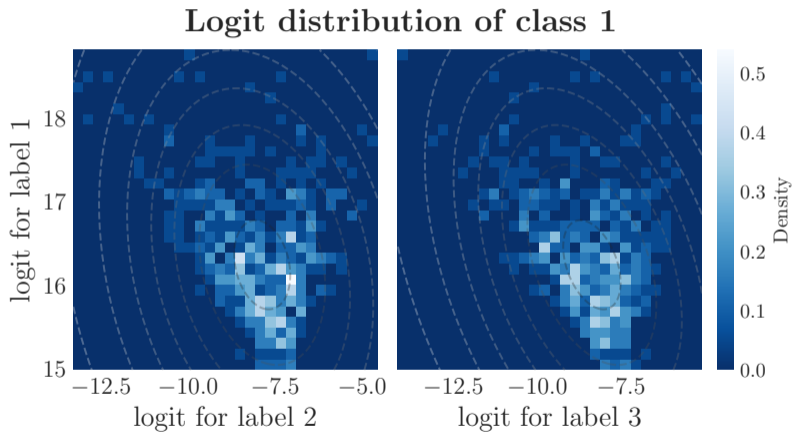


Figure: Joint logit distribution

# Multiclass classification: GMM

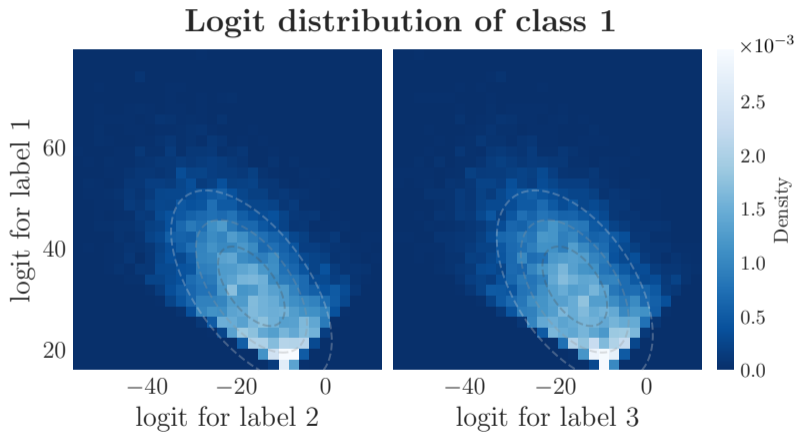


Figure: Joint logit distribution

Thank you for listening.



ArXiv paper



GitHub page

# References

- El Karoui, Nouredine, et al. "On robust regression with high-dimensional predictors." Proceedings of the National Academy of Sciences 110.36 (2013): 14557-14562.
- Donoho, David, and Andrea Montanari. "High dimensional robust m-estimation: Asymptotic variance via approximate message passing." Probability Theory and Related Fields 166 (2016): 935-969.
- Sur, Pragma, and Emmanuel J. Candès. "A modern maximum-likelihood theory for high-dimensional logistic regression." Proceedings of the National Academy of Sciences 116.29 (2019): 14516-14525.
- Belkin, Mikhail, et al. "Reconciling modern machine-learning practice and the classical bias-variance trade-off." Proceedings of the National Academy of Sciences 116.32 (2019): 15849-15854.
- Bartlett, Peter L., et al. "Benign overfitting in linear regression." Proceedings of the National Academy of Sciences 117.48 (2020): 30063-30070.

# References

- Chawla, Nitesh V., et al. "SMOTE: synthetic minority over-sampling technique." Journal of artificial intelligence research 16 (2002): 321-357.
- Zhang, Hongyi, et al. "mixup: Beyond Empirical Risk Minimization." International Conference on Learning Representations. 2018.
- Cao, Kaidi, et al. "Learning imbalanced datasets with label-distribution-aware margin loss." Advances in neural information processing systems 32 (2019).
- Montanari, Andrea, and Kangjie Zhou. "Overparametrized linear dimensionality reductions: From projection pursuit to two-layer neural networks." arXiv preprint arXiv:2206.06526 (2022).