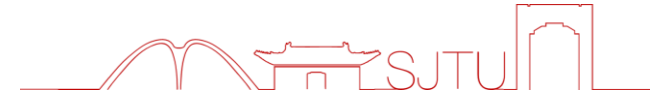




上海交通大学  
SHANGHAI JIAO TONG UNIVERSITY



# Trustworthy Machine Learning on Imbalance Data

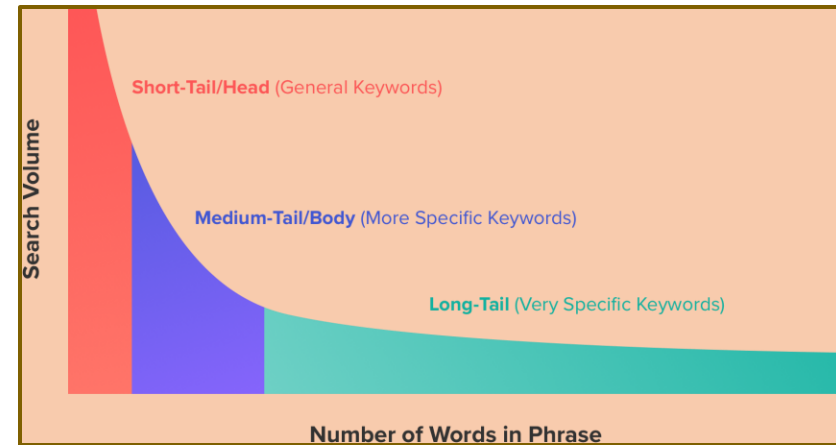
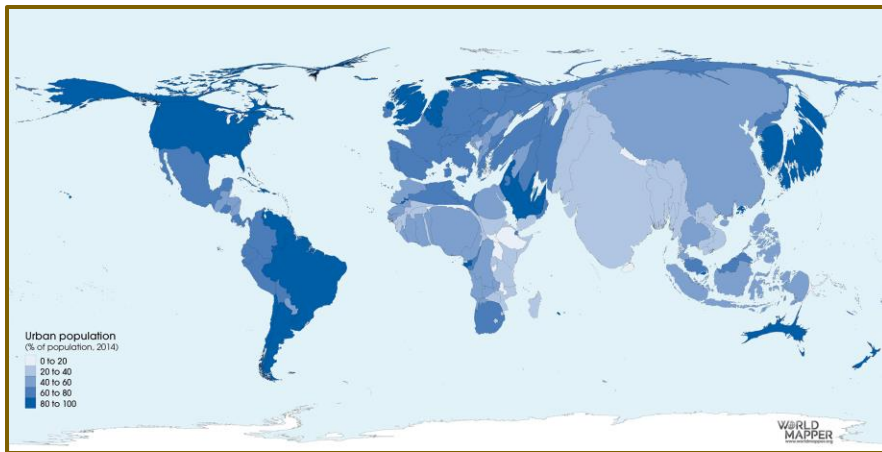
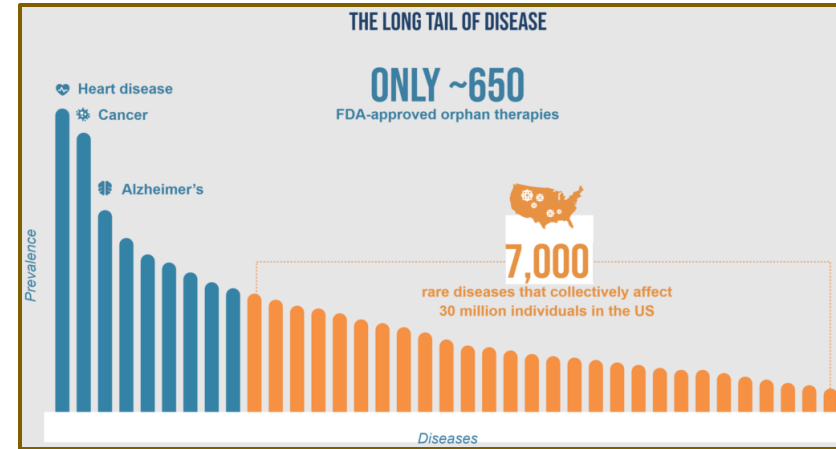
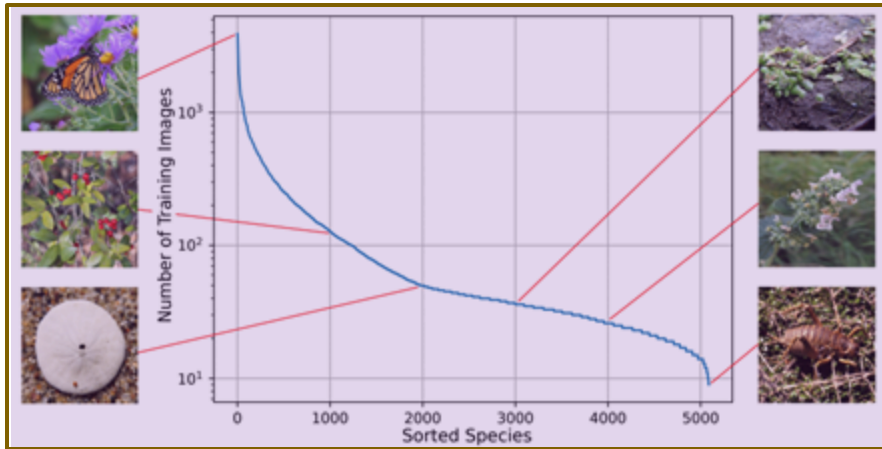
**Jiangchao Yao**

**CMIC, Shanghai Jiao Tong University**

饮水思源 · 爱国荣校



Large-scale natural sources are very imbalance, usually following a **long-tailed distribution**.



[1] Van Horn et al. The iNaturalist Species Classification and Detection Dataset. CVPR 2018.

[2] Gregory et al. CXR-LT challenge. ICCV CVAMD 2023.

[3] <https://worldmapper.org/maps/urban-population-relative-2014/>

[4] <https://seopressor.com/blog/short-tail-or-long-tail-keywords/>



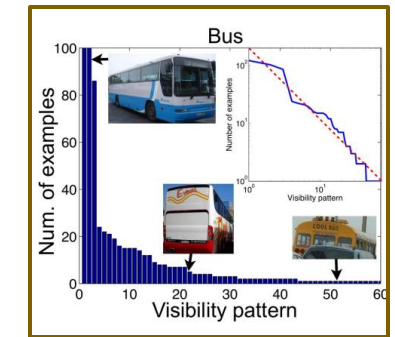
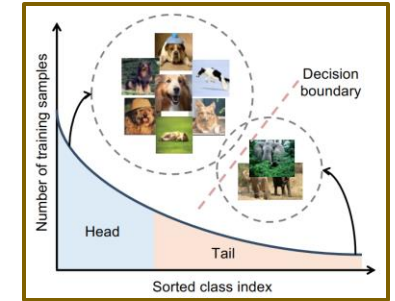
## A direct decomposition on the risk minimization

$$\min_f R(f) = E_{P(x,y)}[\ell(f(x), y)] = \sum_{k=1}^K P(y = k) E_{P(x|y=k)}[\ell(f(x), y)]$$

Minority (“generalized” conceptual) classes have weak importance for training, which can be easily ignored in the early phase especially for overparameterized DNNs [1] (*or namely, will be sacrificed first if it is not sufficient for the model to learn*).

**However**, in real applications, the value of classes cannot be absolutely characterized by their quantity, and instead, sometimes **less is more** for sustainable long-term development.

- **Fairness** w.r.t. diversity e.g., small populations of gender, race and consumers
- **Cost-sensitive scenarios** e.g., medical disease diagnosis and treatment



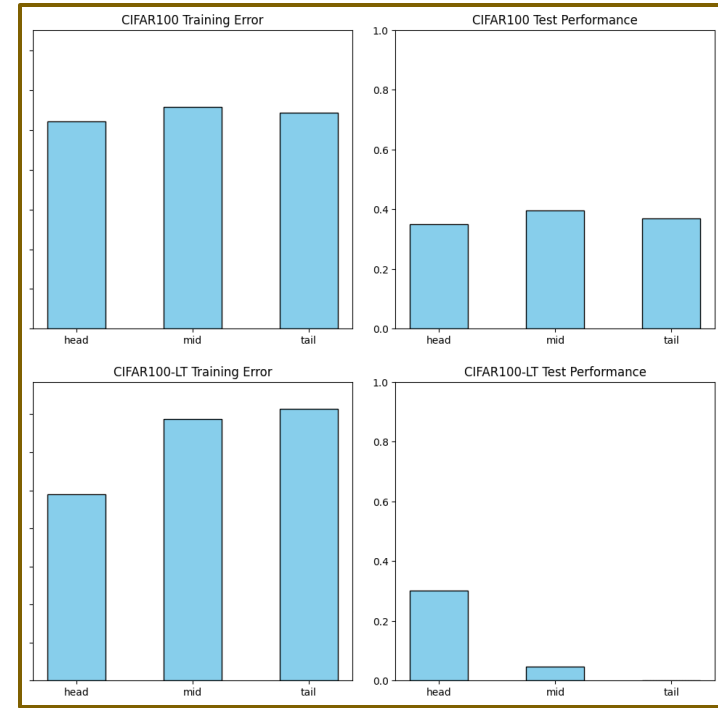
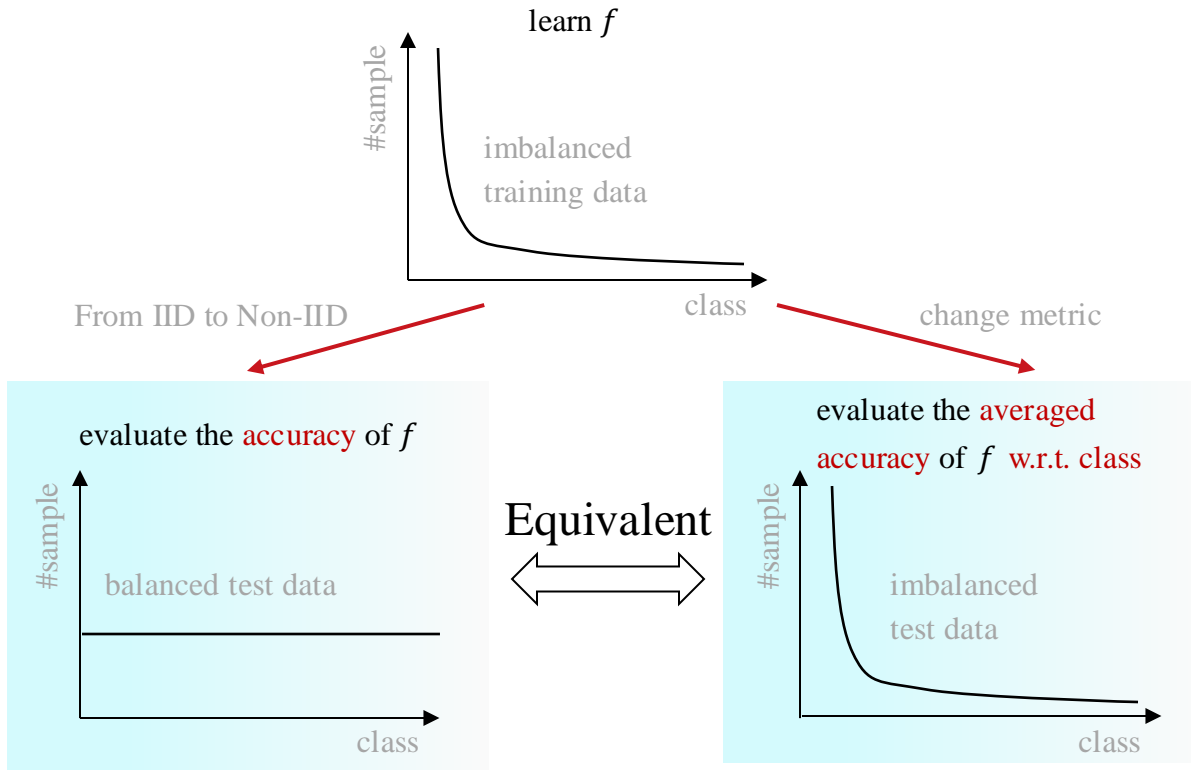
[1] Vitaly Feldman. "Does Learning Require Memorization? A Short Tale about a Long Tail." SIGACT 2020.



# Why the Resulted Imbalanced Learning is Special



A critical highlight on the evaluation, different from the ordinary IID learning



The change in evaluation metric induces **an statistical consistency problem** on applying conventional learning methods, that is,

*What we design during training should be statistically consistent with what we pursue about the evaluation.*





# The Historical Development



## The development of imbalance learning

1998

2009

2018

2023

### Optimizing Classifiers for Imbalanced Training Sets

Grigoris Karakoulas  
Global Analytics Group  
Canadian Imperial Bank of Commerce  
141 Bay St., B.C.E. 11,  
Toronto, ON, Canada M5Z 2S8  
Email: karakou@ci.com

John Shawe-Taylor  
Department of Computer Science  
Royal Holloway, University of London  
Egham, TW20 0EX  
England  
Email: j.shawe@rhul.ac.uk

#### Abstract

Following recent results [8, 8] showing the importance of the fat-shattering dimension in explaining the theoretical effect of a large margin on generalization performance, the current paper investigates the implications of these results for the case of imbalanced datasets and develops two approaches to setting the threshold. The approaches are incorporated into TheRBost, a boosting algorithm for dealing with unequal loss functions. The performance of TheRBost and the two approaches are tested experimentally.

**Keywords:** Computational Learning Theory, Generalization, fat-shattering, large margin, per estimation, unequal loss, imbalanced datasets

#### 1 Introduction

Shawe-Taylor [8] demonstrated that the output margin can also be used as an estimate of the confidence with which a particular classification is made. In other words if a new example has an output value well clear of the threshold we can be more confident of the associated classification than when the output value is close to the threshold. The current paper applies this result to the case where there are different losses associated with a false positive, than with a false negative. If a significant number of data points are misclassified we can use the criterion of minimizing the empirical loss. If, however, the data is correctly classified the empirical loss is zero for all correctly classifying hyperplanes. It is in this case that the approach can provide insight into how to choose the hyperplane and threshold. In summary, the paper suggests ways in which a hyperplane should be optimized for imbalanced datasets where the loss associated with misclassifying the low prevalent class is higher.

### Learning from Imbalanced Data

Habib H. Hamer, IEEE, and Eduardo A. Garcia

**Abstract**—The common occurrence of data available in many scientific, medical, and industrial contexts, such as healthcare, security, financial, or forensic, is heavily skewed in terms of the relative proportion of the two classes. This leads to a significant performance drop in machine learning models trained on such data. This paper reviews the state-of-the-art in this field, focusing on the most recent developments. The paper is organized into two main parts: the first part discusses the theoretical foundations of learning from imbalanced data, and the second part discusses the practical aspects of learning from imbalanced data. The paper is intended as a survey for researchers and practitioners in the field of machine learning, and as a reference for students and researchers in the field of machine learning.

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

#### 1 INTRODUCTION

Recent developments in science and technology have provided the growth and availability of new data sets at an exponential rate. The use of machine learning algorithms to analyze and process this data has led to significant advances in many fields, including healthcare, finance, and transportation. However, the availability of data is often skewed, with one class being much more prevalent than the other. This is known as an imbalanced learning problem, and it poses a significant challenge for machine learning algorithms. In this paper, we review the state-of-the-art in this field, focusing on the most recent developments. The paper is organized into two main parts: the first part discusses the theoretical foundations of learning from imbalanced data, and the second part discusses the practical aspects of learning from imbalanced data. The paper is intended as a survey for researchers and practitioners in the field of machine learning, and as a reference for students and researchers in the field of machine learning.

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

**Index Terms**—Imbalanced learning, distribution, synthetic, methods, cost-sensitive learning, cost-sensitive learning, class imbalance, minority class

Applications

Deep learning

Alberto Fernández · Salvador García  
Mikel Galar · Ronaldó C. Prati  
Bartosz Krawczyk · Francisco Herrera

## Learning from Imbalanced Data Sets

### Deep Long-Tailed Learning: A Survey

Yiwei Zhang, Bingqiang Kong, Binyan Huo, Shucheng Yan, Fuhou IEEE, and Jiahui Peng

**Abstract**—Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

#### 1 Introduction

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

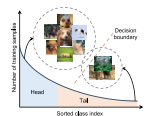


Fig. 1. The distribution of data points in the head and tail classes. The head class has a large number of data points, while the tail class has a small number of data points.

As shown in Fig. 1, we present a typical example of a long-tailed distribution. The head class has a large number of data points, while the tail class has a small number of data points. This is a common scenario in many real-world applications.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

Deep learning has achieved remarkable success in many computer vision tasks, but it still faces the challenge of long-tailed distribution. In this survey, we provide a comprehensive overview of the state-of-the-art in deep long-tailed learning. We first introduce the concept of long-tailed learning and its applications. Then, we review the existing methods for long-tailed learning, including data-level, model-level, and loss-level methods. Finally, we discuss the future research directions in this field.

- [1] Karakoulas et al. Optimizing Classifiers for Imbalanced Training Sets. NIPS 1998.
- [2] He et al. Learning from Imbalanced Data. TKDE 2009.
- [3] Fernández et al. Learning from Imbalanced Data Sets. Springer, 2018.
- [4] Zhang et al. Deep Long-tailed Learning: A Survey. TPAMI 2023.

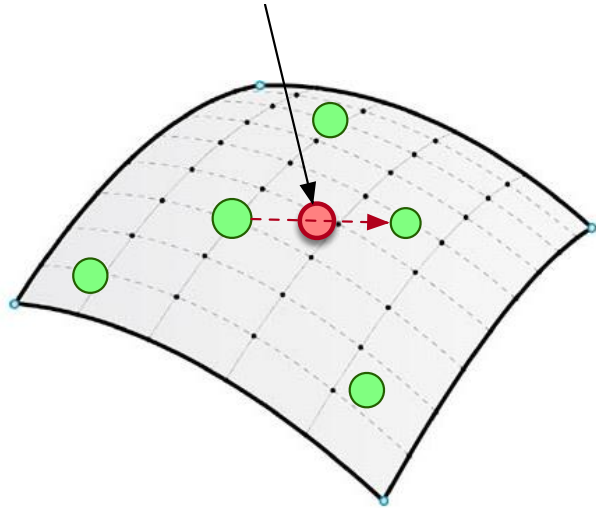




## SMOTE: Synthetic Minority Over-sampling Technique

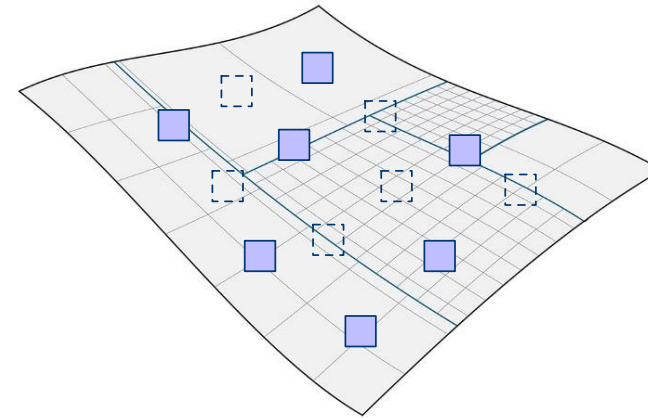
Motivation: Replication of the minority class does not cause its decision boundary to spread into the majority class region (but overfitting).

**Interpolation on minority manifold**



The main idea of SMOTE: augmentation for minority class by interpolation instead of over-sampling with replacement.

**Under-sampling in the majority class**



Interpolation is limited by the samples. Thus, SMOTE also always runs with the under-sampling for majority class.

Increase minority diversity and decrease majority diversity





## Threshold-Moving: adjust the prediction in a post-hoc manner.

Motivation: The over-confident prediction for majority or the low-confident prediction for minority can be calibrated after training.

### THE THRESHOLD-MOVING ALGORITHM

#### Training phase:

1. Let  $S$  be the original training set.
2. Train a neural network from  $S$ .

#### Test phase:

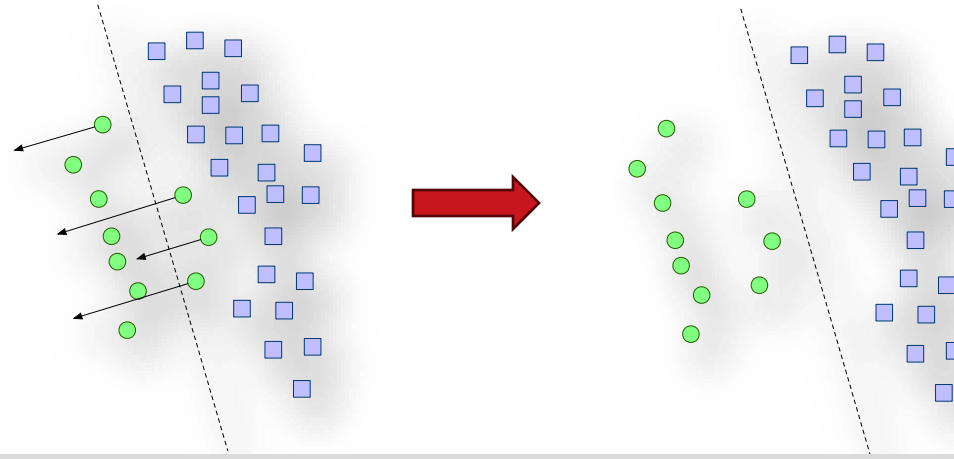
1. Generate real-value outputs with the trained neural network.
2. For every output, multiply it with the sum of the costs of misclassifying the corresponding class to other classes.
3. Return the class with the biggest output.

### Moving function

$$\hat{p}_k = \frac{p_k * \sum_{k'=1}^K C[k][k']}{\eta}$$

where  $p_k$  is the probabilistic prediction,  $C[k][k']$  is the cost mis-predicted from class  $k$  to  $k'$ , and  $\eta$  is renormalization parameter.

Sampling methods might not always show promise in multi-class imbalance learning, but threshold-moving way does.



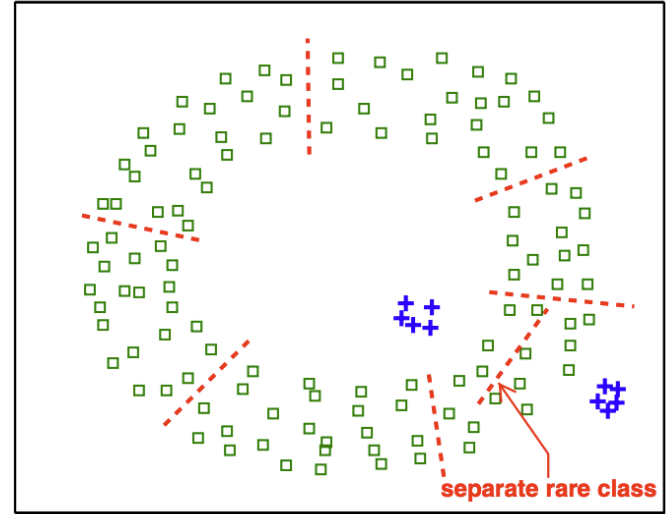
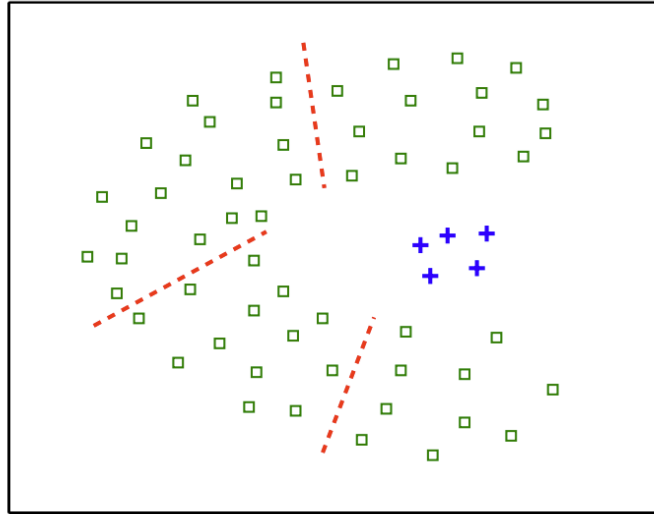
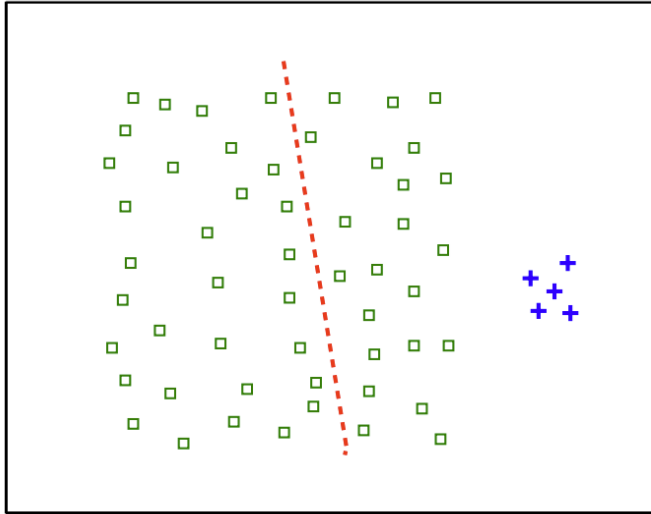
Let  $p_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}$  denote the class prediction. If we set  $\sum_{k'=1}^K C[k][k'] = e^{-\tau \log \pi_k}$  where  $\pi_k$  is the class prior and  $\tau$  is the temperature, the threshold-moving method recovers the popular **logit adjustment** method for long-tailed learning.

Majority classes have the smaller cost than minority classes, e.g.,  $e^{-\tau \log \pi_k}$  is monotonously decreasing.

## COG: Local Decomposition for Rare Class Analysis

Intuition: Quantity imbalance limits the learning pace of minority over majority. We can adjust the quantities by decomposition.

**How to properly decompose the majority classes (or including minority classes) into subclasses to balance the training?**



*Phase I: local clustering*

1. for class  $i = 1$  to  $c$  // " $c$ " represents #classes
2.     clusterLabel( $i$ ) = Clustering( $\mathcal{D}(i)$ ,  $\mathbf{K}(i)$ );
3.      $\mathcal{D}(i)^*$  = changeLabel( $\mathcal{D}(i)$ , clusterLabel( $i$ ));
4. end for

*Phase II: over-sampling (for COG-OS only)*

5. for class  $j = 1$  to  $c$
6.      $\mathcal{D}(j)^{**}$  = replicate( $\mathcal{D}(j)^*$ ,  $\mathbf{r}(j)$ )
7. end for
8.  $\mathcal{D}^{**}$  =  $\bigcup_{j=1}^c (\mathcal{D}(j)^{**})$ ;

*Phase III: training*

9.  $\mathbb{M}$  = train( $\mathcal{D}^{**}$ ,  $\mathbb{L}$ );

*Phase IV: predicting*

10.  $\mathbf{p}'$  = predict( $\mathcal{T}$ ,  $\mathbb{M}$ );
11.  $\mathbf{p}$  = convertLabel( $\mathbf{p}'$ );





## On Statistical Consistency of Binary Classification with Balanced Accuracy

Motivation: The early ERM theory is developed for the instance-wise evaluation, but cannot guarantee the consistency for balanced measure.

$$\text{Accuracy} = \mathbf{E}_{p(x,y)}[h(x) = y]$$



$$\text{Balanced Accuracy} = \frac{\sum_{k \in \{-1,1\}} \mathbf{E}_{p(x|y=k)}[h(x)=k]}{2}$$

If we consider the balanced accuracy, how to modify the algorithm to satisfy the statistical consistency?

**Theorem 3.** Let  $D$  be a probability distribution on  $\mathcal{X} \times \{\pm 1\}$  satisfying Assumption A. Let  $\hat{p}_S$  denote any estimator of  $p = \mathbf{P}(y = 1)$  satisfying  $\hat{p}_S \in (0, 1)$  and  $\hat{p}_S \xrightarrow{P} p$ . Let  $\hat{\eta}_S : \mathcal{X} \rightarrow [0, 1]$  denote any class probability estimator satisfying  $\mathbf{E}_x[|\hat{\eta}_S(x) - \eta(x)|^r] \xrightarrow{P} 0$  for some  $r \geq 1$ , and let  $h_S(x) = \text{sign}(\hat{\eta}_S(x) - \hat{p}_S)$ . Then

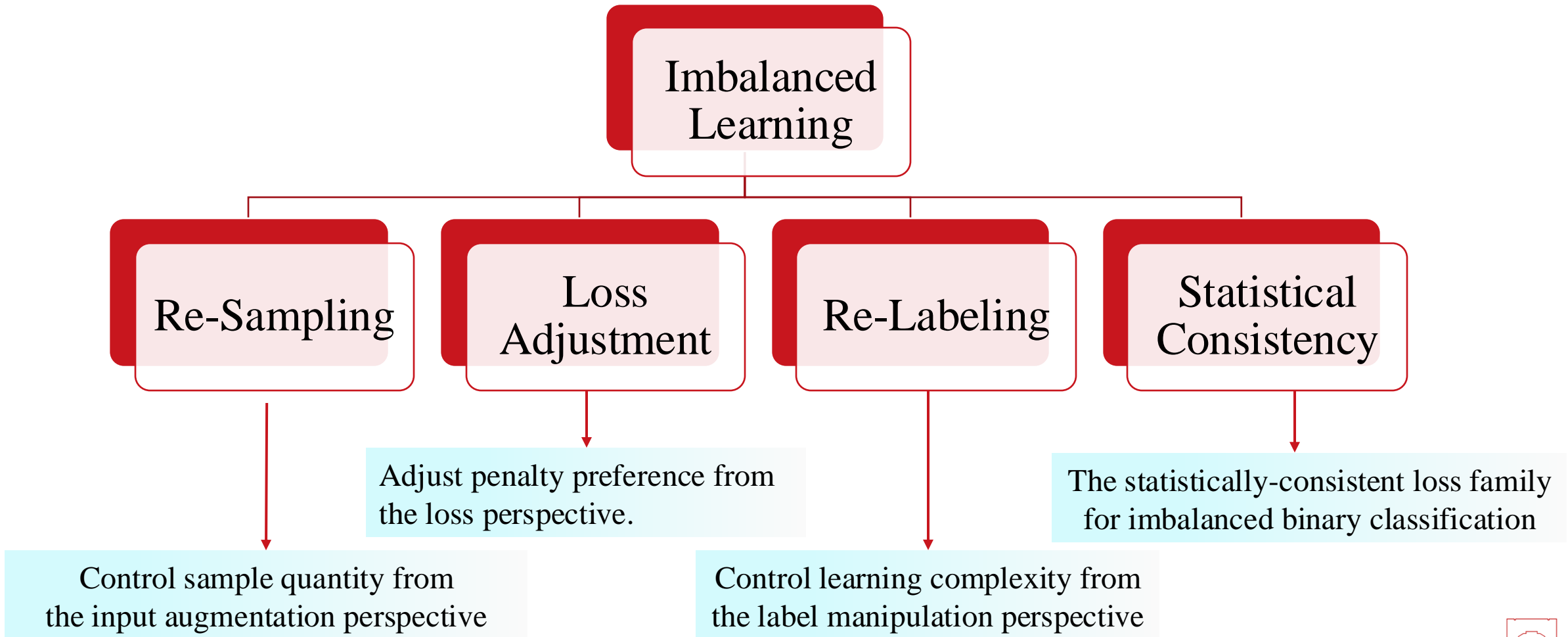
$$\text{regret}_D^{\text{AM}}[h_S] \xrightarrow{P} 0.$$

### Algorithm 1 Plug-in with Empirical Threshold

- 1: **Input:**  $S = ((x_1, y_1), \dots, (x_n, y_n)) \in (\mathcal{X} \times \{\pm 1\})^n$
- 2: **Select:** (a) Proper (composite) loss  $\ell : \{\pm 1\} \times \bar{\mathbb{R}} \rightarrow \bar{\mathbb{R}}_+$ , with link function  $\psi : [0, 1] \rightarrow \bar{\mathbb{R}}$ ; (b) RKHS  $\mathcal{F}_K$  with positive definite kernel  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ ; (c) regularization parameter  $\lambda_n > 0$
- 3:  $f_S \in \text{argmin}_{f \in \mathcal{F}_K} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) + \lambda_n \|f\|_K^2 \right\}$
- 4:  $\hat{\eta}_S = \psi^{-1} \circ f_S$
- 5:  $\hat{p}_S = (\text{as in Eq. (2)})$
- 6: **Output:** Classifier  $h_S(x) = \text{sign}(\hat{\eta}_S(x) - \hat{p}_S)$



## Summary of imbalanced learning in the early years





**What is the new of this topic in the recent years?**

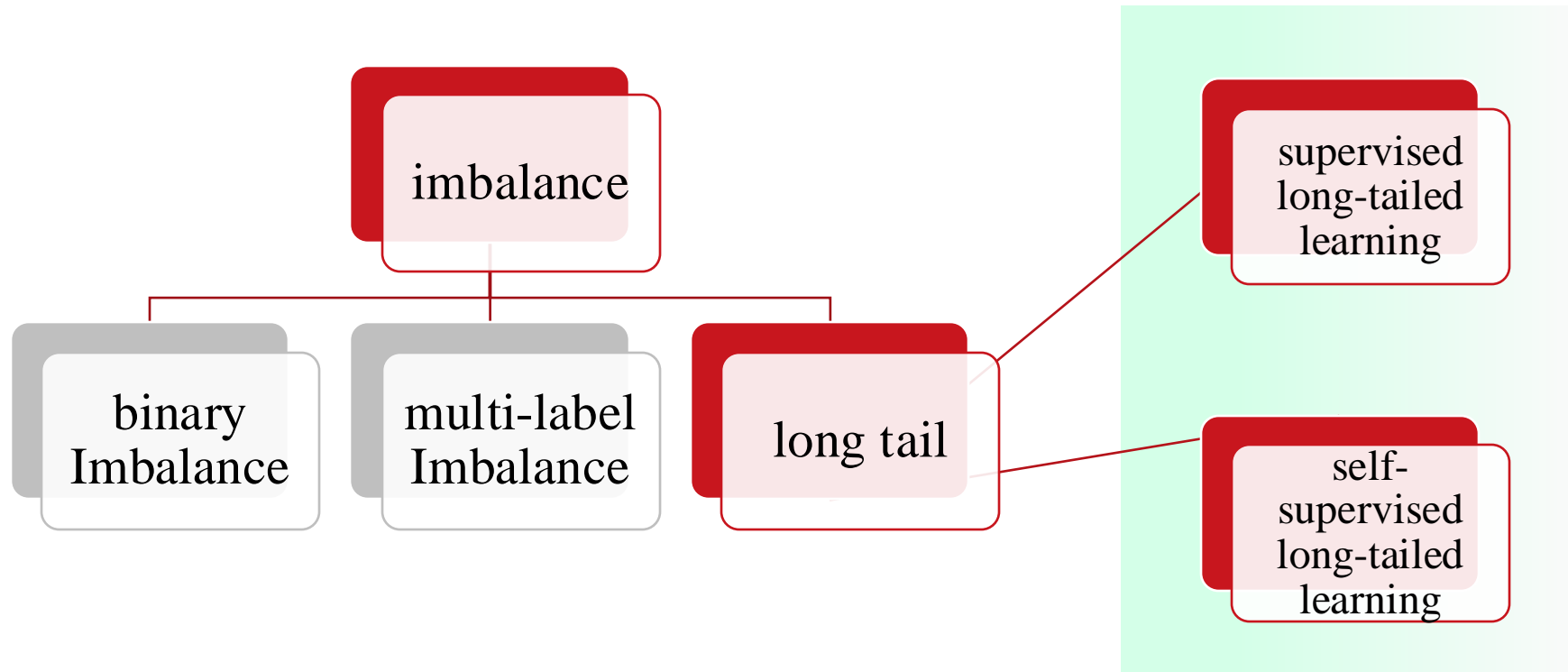




# The Following Part in This Tutorial



The recent advances of imbalance learning powered by deep learning

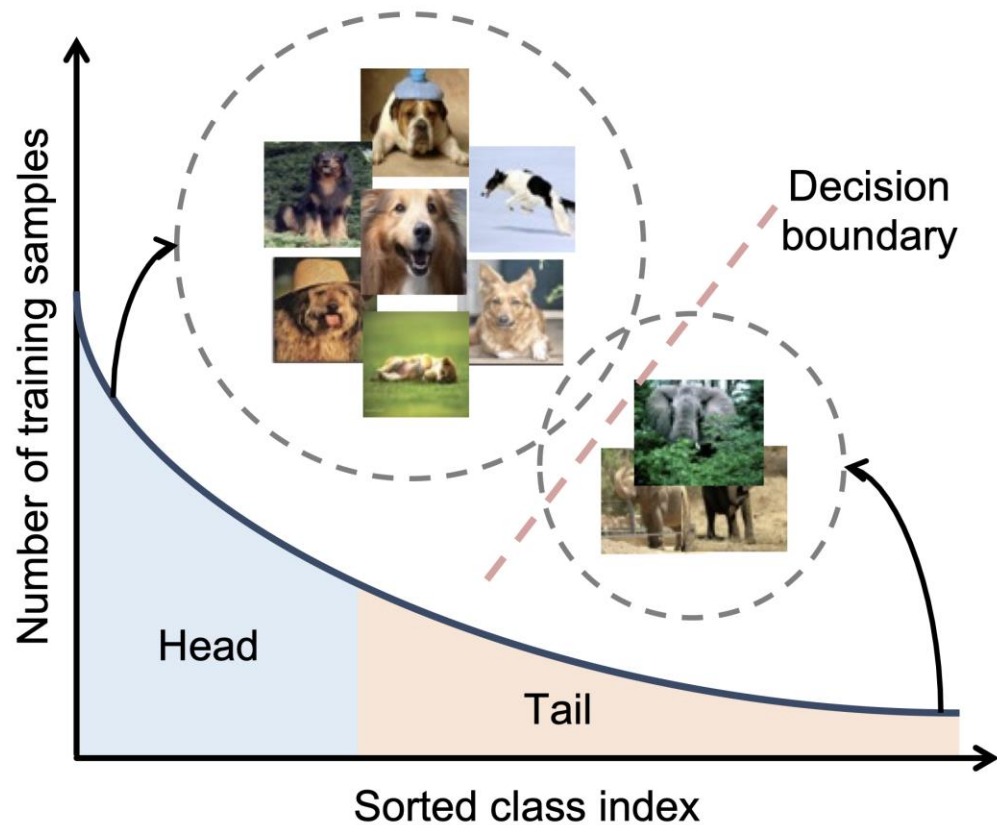




# Supervised Long-tailed Learning



It has been contributed with very broad explorations



Method	Year	Class Re-balancing			Augmentation		Module Improvement				Target Aspect
		Re-sampling	CSL	LA	TL	Aug	RL	CD	DT	Ensemble	
LMLE [89]	2016						✓				feature
HFL [90]	2016						✓				feature
Focal loss [54]	2017		✓								objective
Range loss [21]	2017						✓				feature
CRL [50]	2017						✓				feature
MetaModelNet [91]	2017				✓						
DSTL [92]	2018				✓						
DCL [93]	2019	✓									sample
Meta-Weight-Net [94]	2019		✓								objective
LDAM [18]	2019		✓								objective
CB [16]	2019		✓								objective
UML [95]	2019		✓								feature
FTL [96]	2019		✓		✓	✓					feature
Unequal-training [48]	2019		✓				✓				feature
OLTR [15]	2019		✓				✓				feature
Balanced Meta-Softmax [97]	2020	✓	✓								sample, objective
Decoupling [32]	2020	✓	✓				✓	✓	✓		feature, classifier
LST [98]	2020	✓	✓		✓						sample
Domain adaptation [28]	2020		✓								objective
Equalization loss (ESQL) [19]	2020		✓								objective
DBM [52]	2020		✓								objective
Distribution-balanced loss [37]	2020		✓								objective
UNO-IC [99]	2020		✓								prediction
De-confound-TDE [45]	2020		✓	✓							prediction
M2m [100]	2020		✓		✓						sample
LEAP [49]	2020		✓		✓	✓	✓				feature
OFA [101]	2020		✓		✓	✓			✓		feature
SSP [102]	2020		✓		✓	✓	✓				feature
LFME [103]	2020		✓		✓					✓	sample, model
IEM [104]	2020		✓		✓		✓				feature
Deep-RTC [105]	2020		✓		✓			✓			classifier
SimCal [54]	2020		✓		✓				✓		sample, model
BN [44]	2020		✓		✓				✓		sample, model
BAGS [56]	2020		✓		✓				✓		sample, model
VideoLT [38]	2021		✓		✓						sample
LOCE [33]	2021	✓									sample, objective
DARS [26]	2021	✓	✓		✓						sample, objective
CReST [106]	2021	✓	✓		✓						sample
GIST [107]	2021	✓	✓		✓				✓		classifier
FASA [58]	2021	✓	✓		✓	✓					feature
Equalization loss v2 [108]	2021		✓		✓						objective
Seesaw loss [109]	2021		✓		✓						objective
ACSL [110]	2021		✓		✓						objective
IB [111]	2021		✓		✓						objective
PML [51]	2021		✓		✓						objective
VS [112]	2021		✓		✓						objective
LADE [31]	2021		✓		✓						objective, prediction
RoBal [113]	2021		✓	✓	✓				✓		objective, prediction
DisAlign [29]	2021		✓	✓	✓				✓		objective, classifier
MISLAS [114]	2021		✓	✓	✓				✓		objective, feature, classifier
Logit adjustment [14]	2021		✓	✓	✓						prediction
Conceptual 12M [115]	2021		✓	✓	✓						
DiVE [116]	2021		✓	✓	✓						
MosaicOS [117]	2021		✓	✓	✓						
RSG [118]	2021		✓	✓	✓						feature
SSD [119]	2021		✓	✓	✓						
RIDE [17]	2021		✓	✓	✓				✓		model
MetaAug [120]	2021		✓	✓	✓						sample
PaCo [121]	2021		✓	✓	✓						feature
DRO-LT [122]	2021		✓	✓	✓						feature
Unsupervised discovery [35]	2021		✓	✓	✓						feature
Hybrid [123]	2021		✓	✓	✓						feature
KCL [13]	2021		✓	✓	✓						feature
DT2 [61]	2021		✓	✓	✓				✓		feature, classifier
LTML [46]	2021		✓	✓	✓						sample, model
ACE [124]	2021		✓	✓	✓						sample, model
ResLT [125]	2021		✓	✓	✓						sample, model
SADe [30]	2021		✓	✓	✓						objective, model

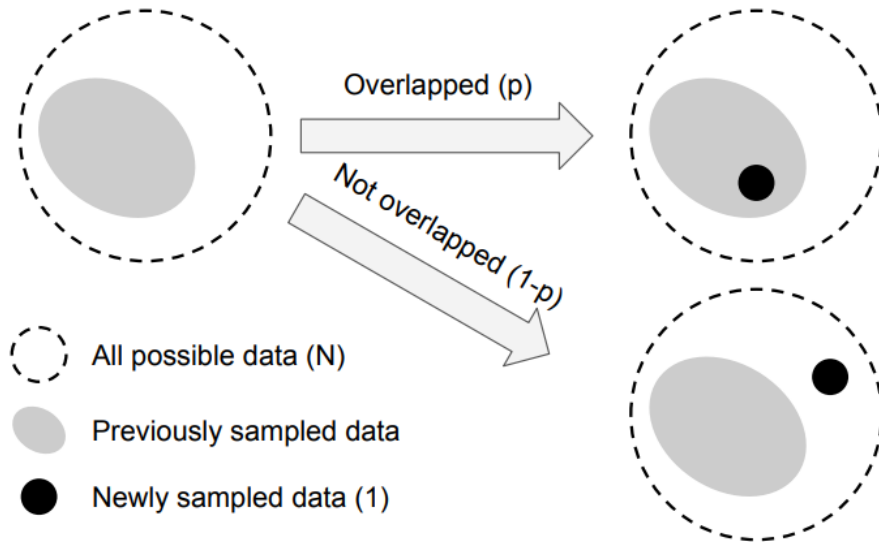
[1] Zhang et al. "Deep Long-Tailed Learning: A Survey" TPAMI 2023.





## loss re-weighting by effective number

➤ **Intuition:** Non-overlapping sample number, instead of the vanilla quantity number, playing the role of imbalance



➤ **Effective Number:** The effective number of examples is the expected volume of samples.

$$E_n = (1 - \beta^n) / (1 - \beta)$$

$$\text{where } \beta = (N - 1) / N$$

$$\lim_{\beta \rightarrow 1} E_n = n$$

➤ **Class-Balanced Loss:** Training from imbalanced data by introducing a weighting factor that is **inversely proportional** to the effective number of samples.

The class-balanced loss term can be applied to a wide range of deep networks and loss functions.





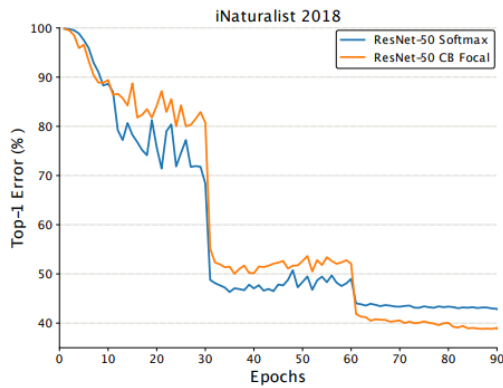
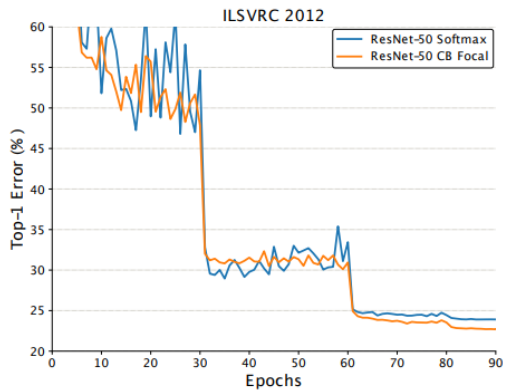
# Supervised Long-tailed Learning



➤ **Class-Balanced Loss:** The class-balanced (CB) loss can be written as:

$$CB(\mathbf{p}, y) = \frac{1}{E_{n_y}} \mathcal{L}(\mathbf{p}, y) = \frac{1 - \beta}{1 - \beta^{n_y}} \mathcal{L}(\mathbf{p}, y) \quad CB_{\text{softmax}}(\mathbf{z}, y) = -\frac{1 - \beta}{1 - \beta^{n_y}} \log \left( \frac{\exp(z_y)}{\sum_{j=1}^C \exp(z_j)} \right)$$

It can also be combined with **sigmoid cross-entropy loss, focal loss**, etc.



Dataset Name	Long-Tailed CIFAR-10						Long-Tailed CIFAR-100					
	Imbalance	200	100	50	20	10	1	200	100	50	20	10
Softmax	<b>34.32</b>	29.64	25.19	17.77	13.61	6.61	65.16	61.68	56.15	48.86	44.29	29.07
Sigmoid	34.51	<b>29.55</b>	23.84	<b>16.40</b>	<b>12.97</b>	<b>6.36</b>	64.39	<b>61.22</b>	55.85	48.57	44.73	<b>28.39</b>
Focal ( $\gamma = 0.5$ )	36.00	29.77	<b>23.28</b>	17.11	13.19	6.75	65.00	61.31	55.88	48.90	44.30	28.55
Focal ( $\gamma = 1.0$ )	34.71	29.62	23.29	17.24	13.34	6.60	<b>64.38</b>	61.59	<b>55.68</b>	<b>48.05</b>	<b>44.22</b>	28.85
Focal ( $\gamma = 2.0$ )	35.12	30.41	23.48	16.77	13.68	6.61	65.25	61.61	56.30	48.98	45.00	28.52
<b>Class-Balanced</b>	<b>31.11</b>	<b>25.43</b>	<b>20.73</b>	<b>15.64</b>	<b>12.51</b>	<b>6.36*</b>	<b>63.77</b>	<b>60.40</b>	<b>54.68</b>	<b>47.41</b>	<b>42.01</b>	<b>28.39*</b>
Loss Type	SM	Focal	Focal	SM	SGM	SGM	Focal	Focal	SGM	Focal	Focal	SGM
$\beta$	0.9999	0.9999	0.9999	0.9999	0.9999	-	0.9	0.9	0.99	0.99	0.999	-
$\gamma$	-	1.0	2.0	-	-	-	1.0	1.0	-	0.5	0.5	-

The proposed framework provides a non-parametric means of quantifying data overlap.



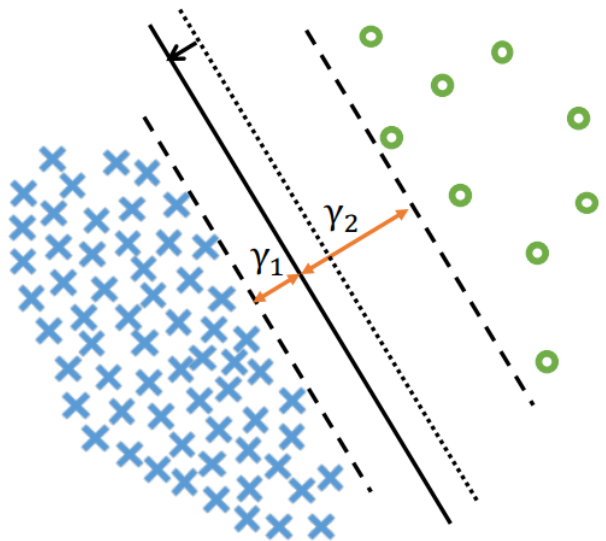


## Class-wise margin calibration

➤ **Motivation:** for imbalanced learning, there is a class-distribution-aware margin trade-off for generalization error.

The generalization error is proportional to the following (two classes, and same holds for multiple classes)

$$\frac{1}{\gamma_1 \sqrt{n_1}} + \frac{1}{\gamma_2 \sqrt{n_2}} \quad \gamma_1 + \gamma_2 = \gamma$$



The margin definition:

$$\gamma_j = \min_{i \in S_j} \gamma(x_i, y_i) \quad \gamma_j = \frac{C}{n_j^{1/4}}$$





# Supervised Long-tailed Learning



➤ **LDAM**: The authors define their hinge loss function (and its relaxed version via Softmax)

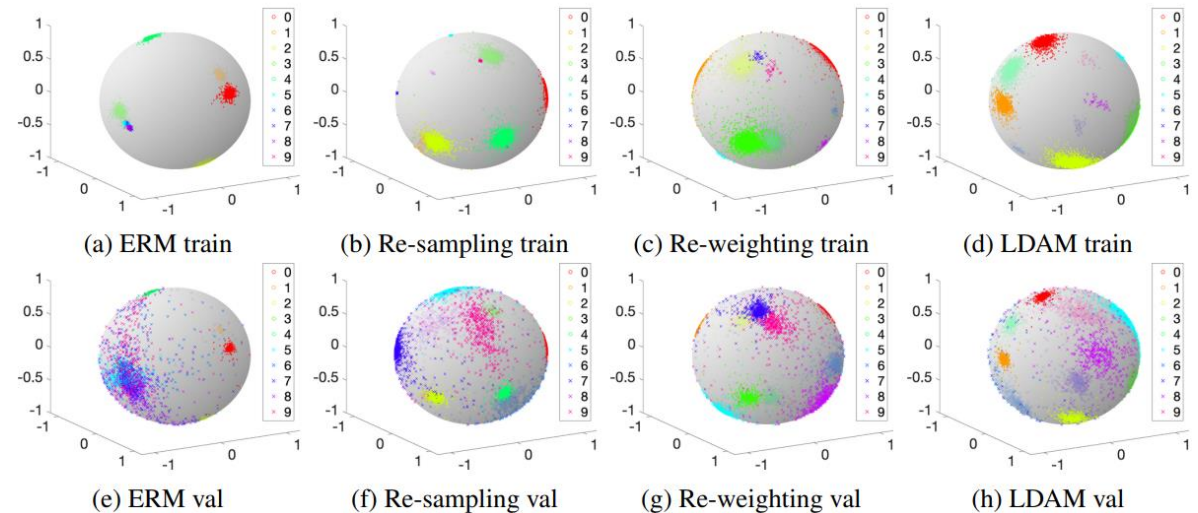
$$\mathcal{L}_{\text{LDAM-HG}}((x, y); f) = \max(\max_{j \neq y} \{z_j\} - z_y + \Delta_y, 0)$$

$$\text{where } \Delta_j = \frac{C}{n_j^{1/4}} \text{ for } j \in \{1, \dots, k\}$$

$$\mathcal{L}_{\text{LDAM}}((x, y); f) = -\log \frac{e^{z_y - \Delta_y}}{e^{z_y - \Delta_y} + \sum_{j \neq y} e^{z_j}}$$

$$\text{where } \Delta_j = \frac{C}{n_j^{1/4}} \text{ for } j \in \{1, \dots, k\}$$

Dataset	Imbalanced CIFAR-10				Imbalanced CIFAR-100			
	long-tailed		step		long-tailed		step	
Imbalance Ratio	100	10	100	10	100	10	100	10
ERM	29.64	13.61	36.70	17.50	61.68	44.30	61.45	45.37
Focal [Lin et al., 2017]	29.62	13.34	36.09	16.36	61.59	44.22	61.43	46.54
LDAM	26.65	13.04	33.42	15.00	60.40	43.09	60.42	43.73
CB RS	29.45	13.21	38.14	15.41	66.56	44.94	66.23	46.92
CB RW [Cui et al., 2019]	27.63	13.46	38.06	16.20	66.01	42.88	78.69	47.52
CB Focal [Cui et al., 2019]	25.43	12.90	39.73	16.54	63.98	42.01	80.24	49.98
HG-DRS	27.16	14.03	29.93	14.85	-	-	-	-
LDAM-HG-DRS	24.42	12.72	24.53	12.82	-	-	-	-
M-DRW	24.94	13.57	27.67	13.17	59.49	43.78	58.91	44.72
<b>LDAM-DRW</b>	<b>22.97</b>	<b>11.84</b>	<b>23.08</b>	<b>12.19</b>	<b>57.96</b>	<b>41.29</b>	<b>54.64</b>	<b>40.54</b>



The margin definition is an approximation to the truth value, and whether we should directly add on the logit space?





## Class-wise logit adjustment [1]

- **Motivation:** Design a consistent loss function that allows for a relatively **elastic margin** in the logit for head and tail.
- **Balanced error:** Under class imbalance, to measure balanced error:

$$\text{BER}(f) \doteq \frac{1}{L} \sum_{y \in [L]} \mathbb{P}_{x|y} \left( y \notin \operatorname{argmax}_{y' \in \mathcal{Y}} f_{y'}(x) \right)$$

Under Bayes-optimal prediction, if  $\mathbb{P}^{\text{bal}}(y | x) \propto \mathbb{P}(y | x) / \mathbb{P}(y)$

Then

$$\operatorname{argmax}_{y \in [L]} \mathbb{P}^{\text{bal}}(y | x) = \operatorname{argmax}_{y \in [L]} \exp(s_y^*(x)) / \mathbb{P}(y) = \operatorname{argmax}_{y \in [L]} s_y^*(x) - \ln \mathbb{P}(y)$$





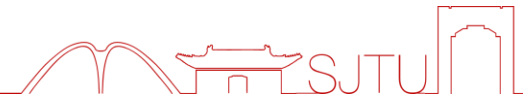
## ➤ The logit adjusted softmax cross-entropy

$$\ell(y, f(x)) = -\log \frac{e^{f_y(x) + \tau \cdot \log \pi_y}}{\sum_{y' \in [L]} e^{f_{y'}(x) + \tau \cdot \log \pi_{y'}} = \log \left[ 1 + \sum_{y' \neq y} \left( \frac{\pi_{y'}}{\pi_y} \right)^\tau \cdot e^{(f_{y'}(x) - f_y(x))} \right]$$

$$w_1^\top \Phi(x) / \pi_1 < w_2^\top \Phi(x) / \pi_2 \not\Rightarrow \exp(w_1^\top \Phi(x)) / \pi_1 < \exp(w_2^\top \Phi(x)) / \pi_2.$$

## ➤ Post-hoc logit adjustment

$$\operatorname{argmax}_{y \in [L]} \exp(w_y^\top \Phi(x)) / \pi_y^\tau = \operatorname{argmax}_{y \in [L]} f_y(x) - \tau \cdot \log \pi_y$$





- An remarkable point on the statistical consistency of long-tailed multi-class classification

$$\ell(y, f(x)) = \alpha_y \cdot \log \left[ 1 + \sum_{y' \neq y} e^{\Delta_{yy'}} \cdot e^{(f_{y'}(x) - f_y(x))} \right]$$



**Theorem 1.** For any  $\delta \in \mathbb{R}_+^L$ , the pairwise loss in (11) is Fisher consistent with weights and margins

$$\alpha_y = \delta_y / \mathbb{P}(y) \quad \Delta_{yy'} = \log(\delta_{y'} / \delta_y).$$

Letting  $\delta_y = \pi_y$ , we immediately deduce that the logit-adjusted loss of (10) is consistent, *provided* our  $\pi_y$  is a consistent estimate of  $\mathbb{P}(y)$ . Similarly,  $\delta_y = 1$  recovers the classic result that the balanced loss is consistent. While Theorem 1 only provides a sufficient condition in multi-class setting, one can provide a necessary and sufficient condition that rules out other choices of  $\Delta$  in the binary case.







## Dynamic adjustment based on a fine-grained generalization bound

**Proposition 3** (Data-Dependent Bound for the VS Loss). *Given the function set  $\mathcal{F}$  and the VS loss  $L_{VS}$ , for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  over the training set  $\mathcal{S}$ , the following generalization bound holds for all  $f \in \mathcal{F}$ :*

$$\mathcal{R}_{bal}^L(f) \lesssim \Phi(L_{VS}, \delta) + \frac{\hat{\mathcal{E}}_{\mathcal{S}}(\mathcal{F})}{C\pi_C} \sum_{y=1}^C \alpha_y \tilde{\beta}_y \sqrt{\pi_y} [1 - \text{softmax}(\beta_y B_y(f) + \Delta_y)].$$

$$L_{VS}(f(\mathbf{x}), y) = -\alpha_y \log \left( \frac{e^{\beta_y f(\mathbf{x})_y + \Delta_y}}{\sum_{y'} e^{\beta_{y'} f(\mathbf{x})_{y'} + \Delta_{y'}}} \right).$$

---

**Algorithm 1:** Principled Learning Algorithm induced by the Theoretical Insights

---

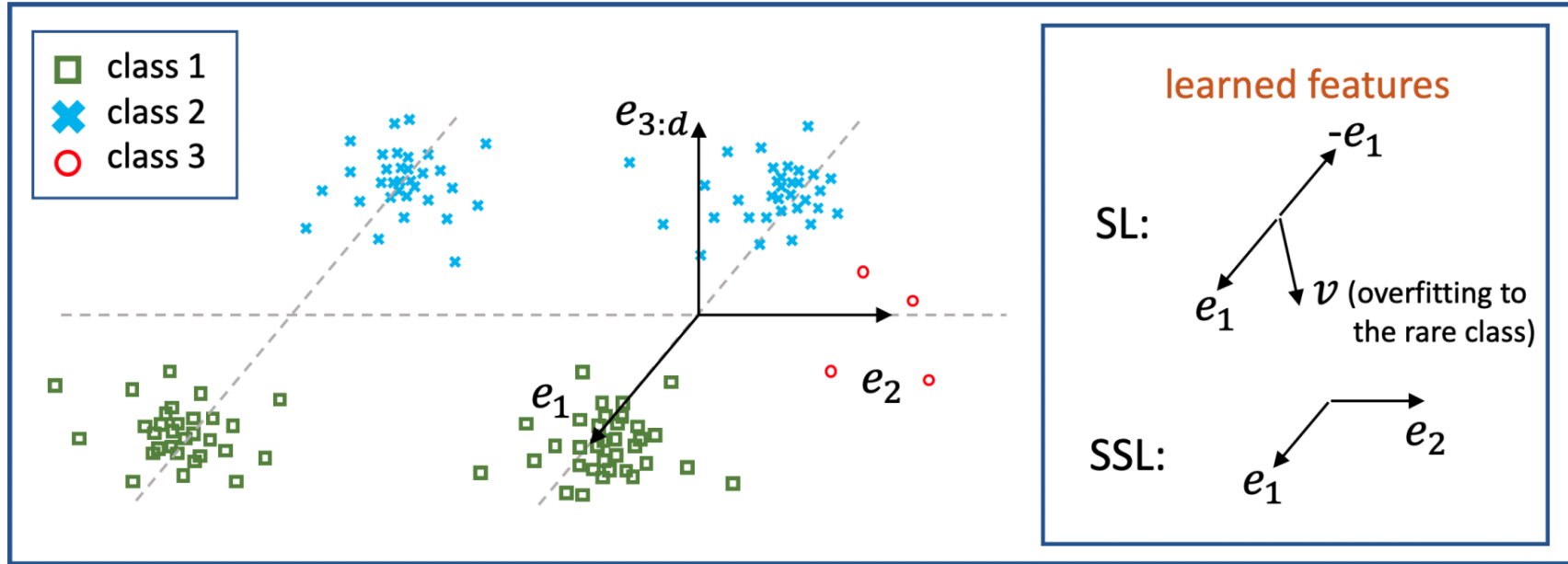
**Require:** Training set  $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^N$  and a model  $f$  parameterized by  $\Theta$ .

- 1: Initialize the model parameters  $\Theta$  randomly.
  - 2: **for**  $t = 1, 2, \dots, T$  **do**
  - 3:    $\mathcal{B} \leftarrow \text{SampleMiniBatch}(\mathcal{S}, m)$  ▷ A mini-batch of  $m$  samples
  - 4:   **if**  $t < T_0$  **then**
  - 5:     Set  $\alpha = 1, \beta_y, \Delta_y$  ▷ Adjust logits during the initial phase
  - 6:   **else**
  - 7:     Set  $\alpha_y \propto \pi_y^{-\nu}, \beta_y = 1, \Delta_y, \nu > 0$  ▷ TLA and ADRW
  - 8:   **end if**
  - 9:    $L(f, \mathcal{B}) \leftarrow \frac{1}{m} \sum_{(\mathbf{x}, y) \in \mathcal{B}} L_{VS}(f(\mathbf{x}), y)$  ▷ Calculate the loss
  - 10:    $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} L(f, \mathcal{B})$  ▷ One SGD step
  - 11:   Optional: anneal the learning rate  $\eta$ . ▷ Required when  $t = T_0$
  - 12: **end for**
- 





Is self-supervised learning more robust to data imbalance?



- Supervised learning (SL) only **extracts features that are useful for predicting labels** ( $e_1$ )

- Self-supervised learning (SSL) learns **task-irrelevant features regardless of the labels**, which enables richer and more robust representation ( $e_1, e_2$ )

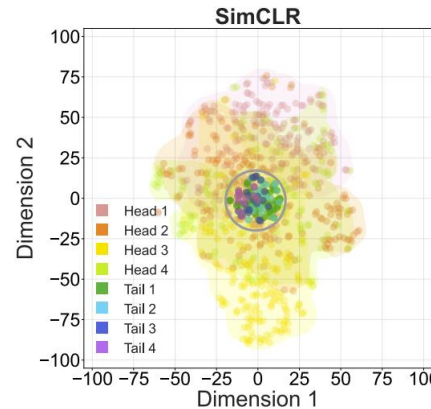
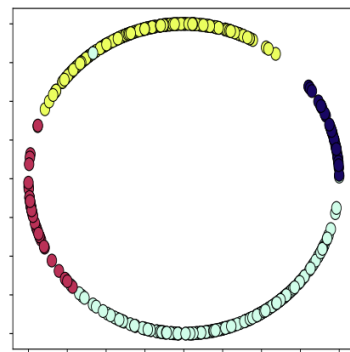
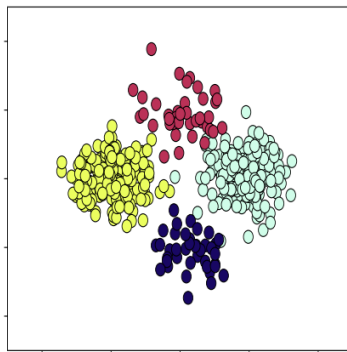


## Self-supervised learning still suffers from data imbalance

➤ Performance degeneration: Linear probing on imbalanced data ( $D_i$ ) and balanced data ( $D_b$ ) with same data amount

Dataset	Subset	<i>Many</i>	<i>Medium</i>	<i>Few</i>	All
CIFAR10	$D_b$	$77.14 \pm 4.64$	$74.25 \pm 6.54$	$71.47 \pm 7.55$	$74.57 \pm 0.65$
	$D_i$	$76.07 \pm 3.88$	$67.97 \pm 5.84$	$54.21 \pm 10.24$	$67.08 \pm 2.15$
CIFAR100	$D_b$	$25.48 \pm 1.74$	$25.16 \pm 3.07$	$24.01 \pm 1.23$	$24.89 \pm 0.99$
	$D_i$	$30.72 \pm 2.01$	$21.93 \pm 2.61$	$15.99 \pm 1.51$	$22.96 \pm 0.43$

➤ Representation learning disparity: head classes dominate the feature regime but tail classes passively collapse



[1] Jiang et al. "Self-damaging contrastive learning." ICML 2021.

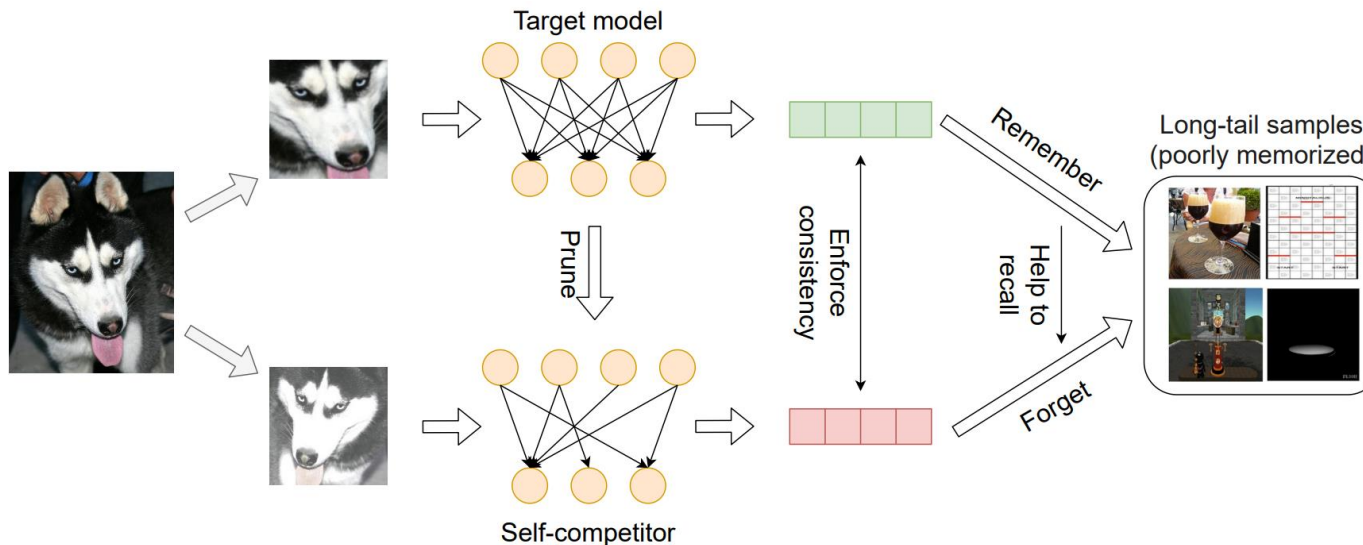
[2] Zhou et al. "Combating Representation Learning Disparity with Geometric Harmonization." NeurIPS 2023.





## SDCLR: Self-damaging Contrastive Learning

- **Intuition:** The sensitivity of head and tail samples to the model pruning, are very different, which helps us to anchor and promote the training of tail samples.
- **Pruning identified exemplars (PIE)** systematically investigates the model output changes introduced by pruning and finds that certain examples are particularly sensitive to sparsity. **They are high likely to be rare and atypical samples, which probably comes from tail classes.**

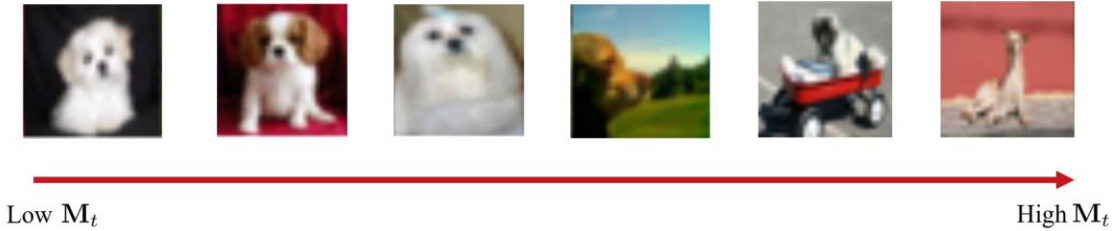
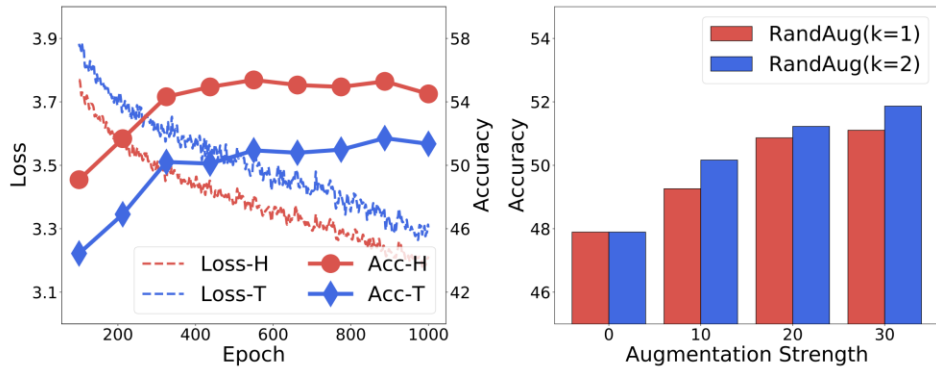




## BCL: Boosted Contrastive Learning

- Motivation I: Memorization effect still holds under long-tailed distribution.
- Motivation II: Stronger information discrepancy motivates tail samples mining

➤ **Challenge:** how to **detect tail data** and how to construct the **desired information discrepancy**



- Motivated from the observation that *learning speed-based proxy* shows strong correlation with the memorization score[1], BCL extends the memorization estimation to *self-supervised learning*.

$$\mathcal{L}_{i,0}^m = \mathcal{L}_{i,0}, \quad \mathcal{L}_{i,t}^m = \beta \mathcal{L}_{i,t-1}^m + (1 - \beta) \mathcal{L}_{i,t}$$

$$\mathbf{M}_{i,t} = \frac{1}{2} \left( \frac{\mathcal{L}_{i,t}^m - \bar{\mathcal{L}}_t^m}{\max \{ |\mathcal{L}_{i,t}^m - \bar{\mathcal{L}}_t^m| \}_{i=0, \dots, N}} + 1 \right)$$

$$\Psi(x_i; \mathcal{A}, \mathbf{M}_i) = a_1(x_i) \circ \dots \circ a_k(x_i),$$

$$a_j(x_i) = \begin{cases} A_j(x_i; \mathbf{M}_i \zeta) & u \sim \mathcal{U}(0, 1) \ \& \ u < \mathbf{M}_i \\ x_i & \text{otherwise} \end{cases}$$

$$\mathcal{L}_{\text{BCL}} = \frac{1}{N} \sum_{i=1}^N -\log \frac{\exp \left( \frac{f(\Psi(x_i))^\top f(\Psi(x_i^+))}{\tau} \right)}{\sum_{x'_i \in X'} \exp \left( \frac{f(\Psi(x_i))^\top f(\Psi(x'_i))}{\tau} \right)}$$

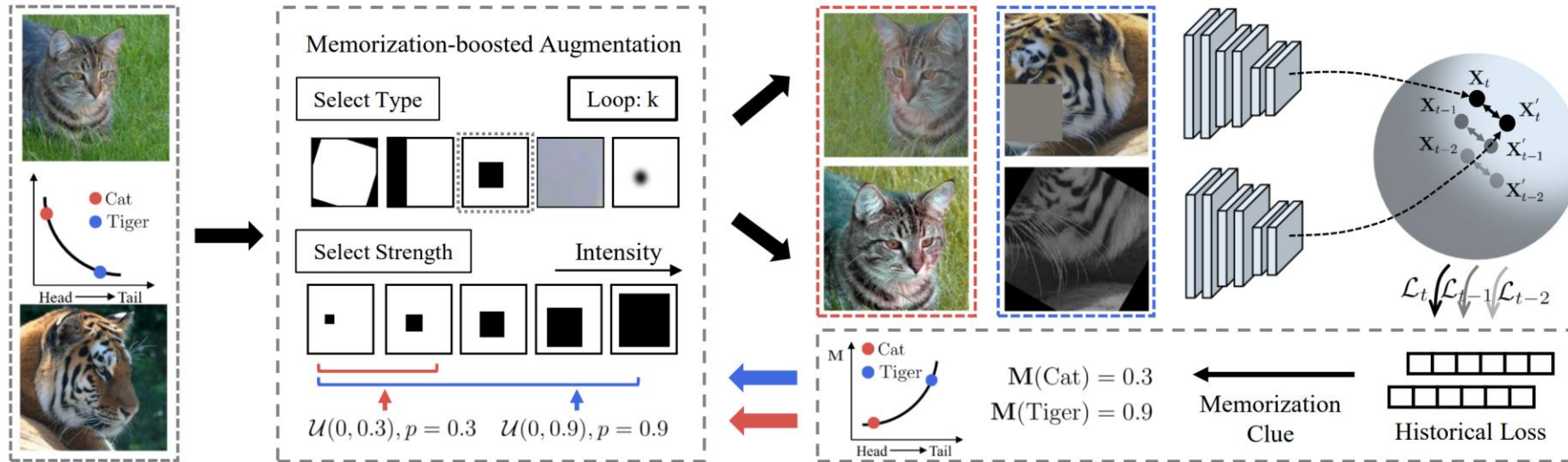
Adaptively assigns the appropriate augmentation strength for the individual sample according to the feedback from the memorization clues

[1] Jiang et al. "Characterizing structural regularities of labeled data in overparameterized models." ICML 2021

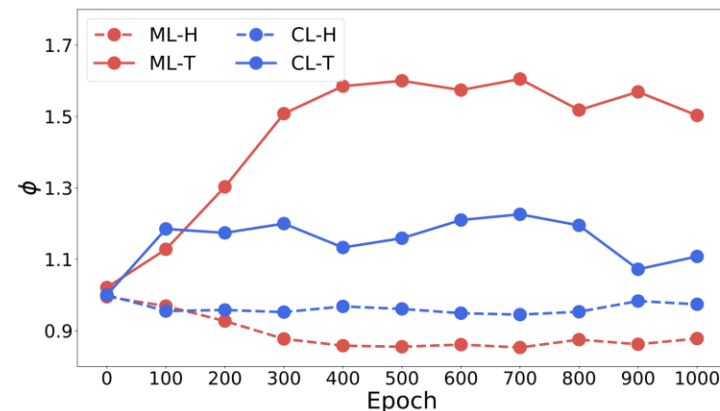
[2] Zhou et al. "Contrastive learning with boosted memorization." ICML 2022.



## BCL: Boosted Contrastive Learning



- Calculate **memorization scores** based on historical statistics to detect tail.
- Construct **instance-wise augmentations** to enhance representation learning.



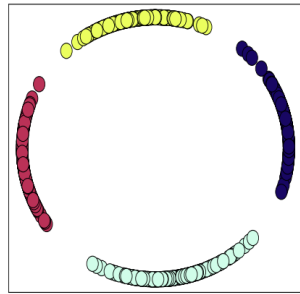
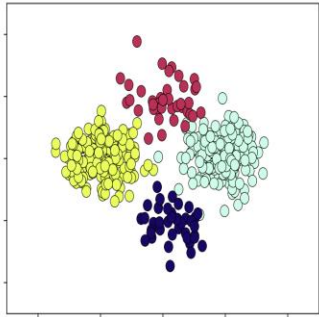




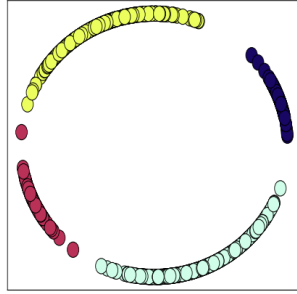
## GH: Geometric Harmonization

➤ Why the conventional contrastive learning underperforms in self-supervised long-tailed context?

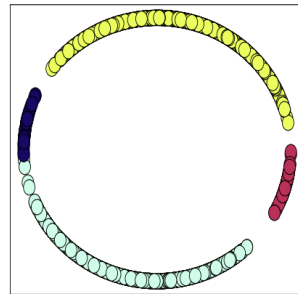
Conventional contrastive loss motivates *sample-level uniformity*, which is biased towards the head classes.



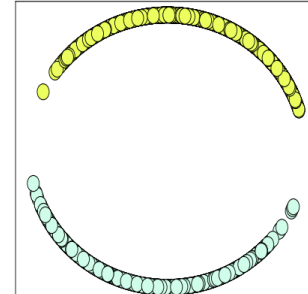
(a) R=1(Balanced)



(b) R=4



(c) R=16



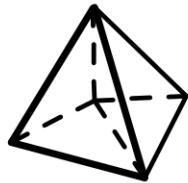
(d) R=64

Contrastive learning causes severer representation learning disparity when enlarging the imbalance ratios.

### Geometric Uniform Structure

$$\mathbf{M}_i^\top \cdot \mathbf{M}_j = C, \quad \forall i, j \in \{1, 2, \dots, K\}, i \neq j,$$

Any two vectors in  $\mathbf{M}$  have the same angle, namely, the unit space are equally partitioned by the vectors.



### Surrogate Label Allocation

$$\min_{\hat{\mathbf{Q}}=[\hat{q}_1, \dots, \hat{q}_N]} \mathcal{L}_{\text{GH}} = -\frac{1}{|\mathcal{D}|} \sum_{x_i \sim \mathcal{D}} \hat{q}_i \log q_i,$$

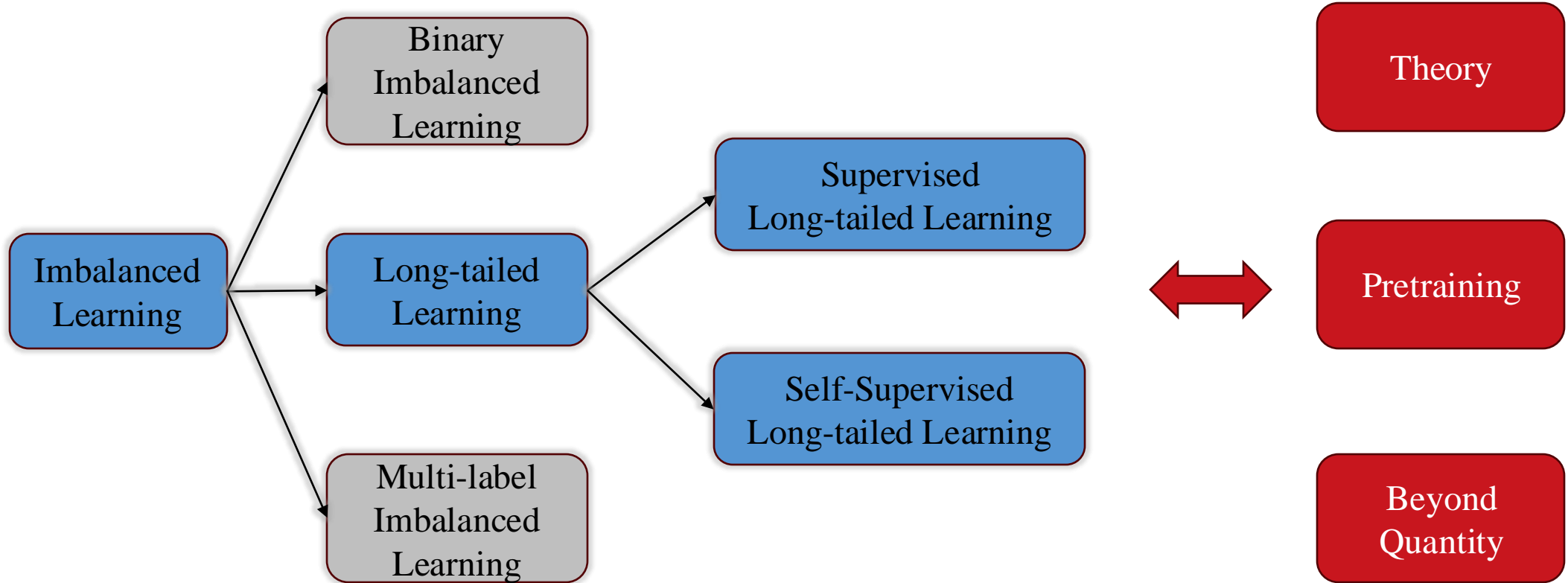
$$\text{s.t. } \hat{\mathbf{Q}} \cdot \mathbf{1}_N = N \cdot \boldsymbol{\pi}, \quad \hat{\mathbf{Q}}^\top \cdot \mathbf{1}_K = \mathbf{1}_N,$$

### Overall objective

$$\min_{\theta, \hat{\mathbf{Q}}} \mathcal{L} = \mathcal{L}_{\text{InfoNCE}} + w_{\text{GH}} \mathcal{L}_{\text{GH}},$$



Still require more efforts on this way





Thank you

Q & A

