

Long-tail learning via logit adjustment

Aditya Krishna Menon



Sadeep Jayasumana



Ankit Singh Rawat



Himanshu Jain



Andreas Veit

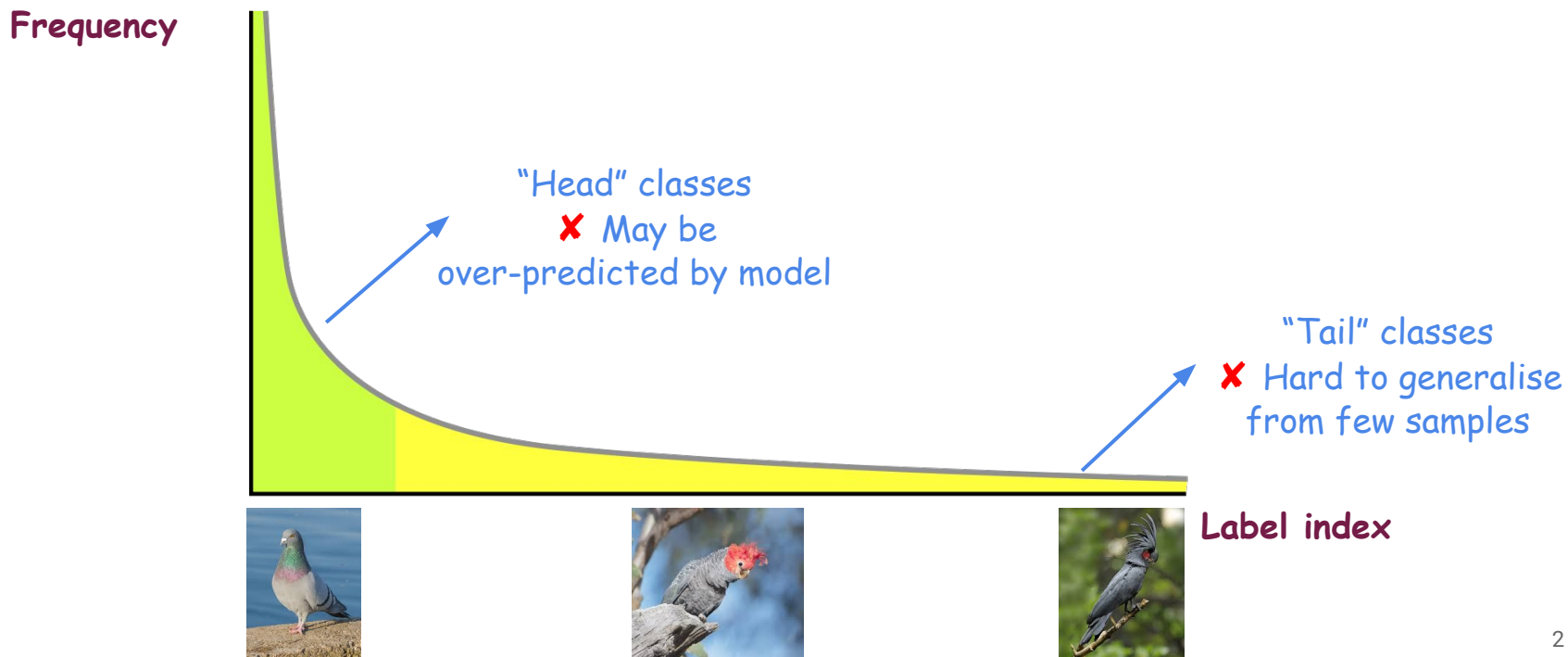


Sanjiv Kumar



Long-tail learning

Classification where the label distribution is **skewed**

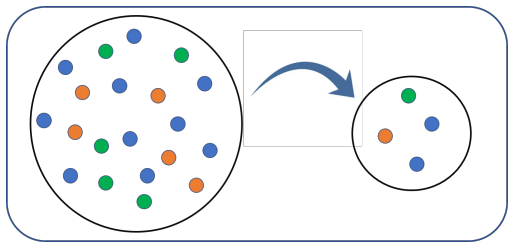


Summary of our work

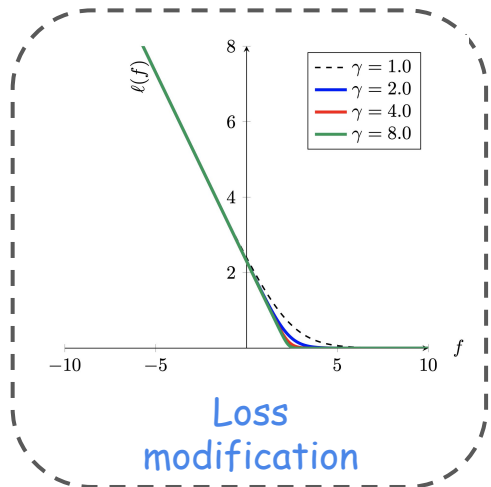
- ✓ A **statistical perspective** of long-tail learning
- ✓ Unifies and **generalises** existing approaches
- ✓ Yields new **post-hoc** and **loss modification** approaches

Existing approaches

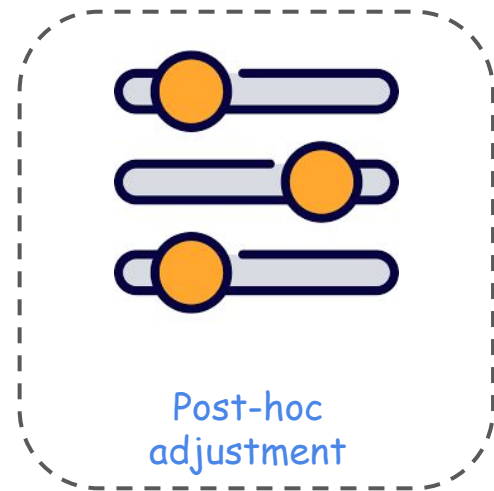
General strategies



Data sampling



Loss modification



Post-hoc adjustment

Weight normalisation

For instance x , a neural model computes logits

$$f_y(x) = w_y^T \Phi(x)$$

✘ Norms may be smaller
for rare classes!

Weight normalisation [Kang et al., '20]:

- (1) Learn w, Φ via standard ERM
- (2) Post-hoc **normalise** the weight norms

Loss modification

Enforce varying margin depending on label frequency:

Softmax:

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

e.g., $1/P(y)$

Adaptive margin:

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

Encourage higher margin between rare +ve and all -ves

[Cao et al., 2019]

e.g., $P(y')$

Equalised loss:

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

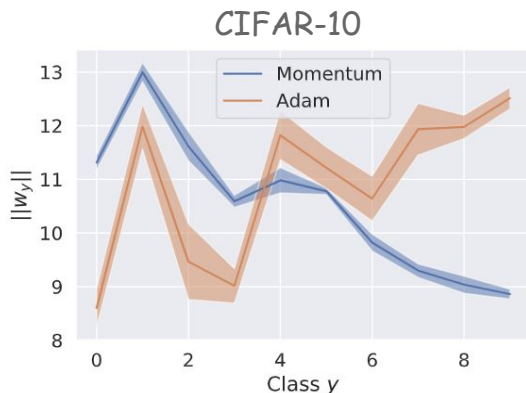
Prevent rare -ves from having gradient overwhelmed

[Tan et al., 2020]
Google

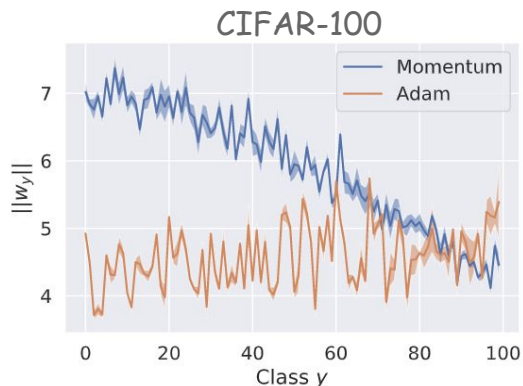
Weight normalisation: limitation

$$f_y(x) = w_y^T \Phi(x)$$

✗ Norms may be smaller for rare classes!



Classes ordered by frequency



Classes ordered by frequency


Weight norms don't correlate with $P(y)$ when using Adam

A statistical framework

Balanced error


Typically, measure **misclassification** error:

$$\begin{aligned}\text{ERR}(h) &= \mathbb{P}(y \neq h(x)) \\ &= \sum_{i \in [L]} \mathbb{P}(y = i) \cdot \mathbb{P}(y \neq h(x) \mid y = i)\end{aligned}$$

 **✗** Can do very well by predicting majority class!

Under class imbalance, can measure **balanced error**:

$$\text{BER}(h) = \sum_{i \in [L]} \frac{1}{L} \cdot \mathbb{P}(y \neq h(x) \mid y = i)$$

 **✓** Treat all classes equally

Statistical view of long-tail learning

Bayes-optimal prediction:

$$\operatorname{argmax}_{y \in [L]} \mathbb{P}^{\text{bal}}(y | x)$$

↓
Balanced label
distribution

Equivalently, if $\mathbf{P}(y | x) \propto \exp(s_y^*(x))$,

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

↓
"Optimal" logits

↓
Increase score for
rare classes

Strategies for long-tail learning

Bayes-optimal solution suggests two strategies:

- (1) Estimate $\mathbf{P}(y | x)$, and adjust logits **post-hoc**
- (2) Directly estimate $\mathbf{P}_{\text{bal}}(y | x)$ by **inherently** adjusting logits

Post-hoc logit adjustment

Standard prediction:

$$\operatorname{argmax}_{y \in [L]} \exp(w_y^T \Phi(x)) = \operatorname{argmax}_{y \in [L]} f_y(x)$$

Logit adjusted prediction:

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

Scaling parameter

Estimate of $P(y)$

When $\tau > 1$, equivalent to [temperature-scaling](#) the probabilities

Comparison to weight normalisation

Logit adjustment performs **additive** correction:

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

Weight normalisation performs **multiplicative** correction:

$$\operatorname{argmax}_{y \in [L]} (w_y^T \Phi(x)) / \pi_y^\tau = \operatorname{argmax}_{y \in [L]} f_y(x) / \pi_y^\tau$$

Logit adjusted loss

Logit adjusted softmax cross-entropy:

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

↙
Add fixed offset to logits

↘
> 1 when $P(y') > P(y)$, i.e., -ve
is more common than +ve

Now predict $\operatorname{argmax}_y f_y(x)$ as normal

A margin view

Consider the **pairwise margin loss**

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

Existing losses $\rightarrow \Delta_{yy}$ depends on y or y' , but not both

Logit adjustment $\rightarrow \Delta_{yy'} = \log \mathbf{P}(y')/\mathbf{P}(y) = \log \mathbf{P}(y') - \log \mathbf{P}(y)$

Enforces a **relative margin** between labels

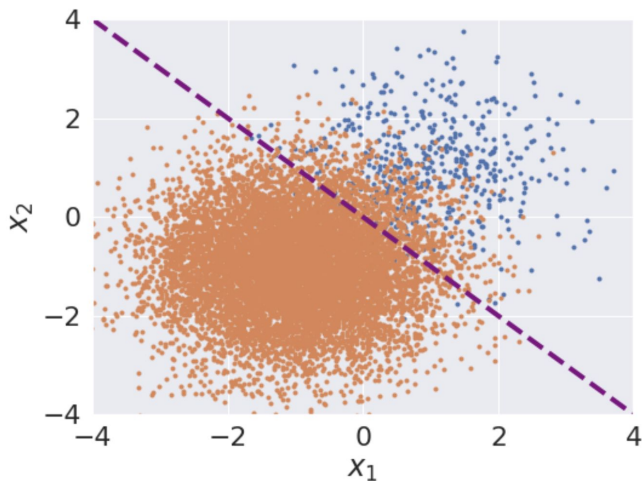
Experiments

Experiments: synthetic data

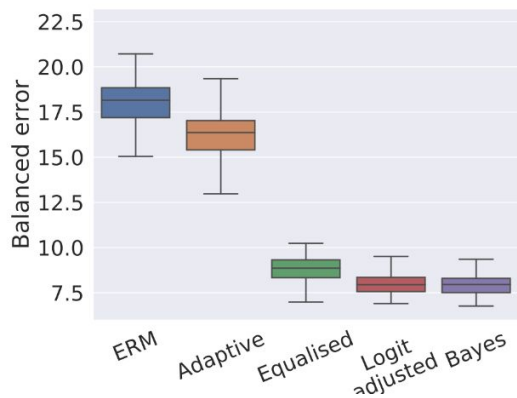
Consider data drawn from a mixture of isotropic Gaussians, with $\mathbf{P}(y = 1) = 5\%$

Bayes-optimal for balanced error: separator passing through origin

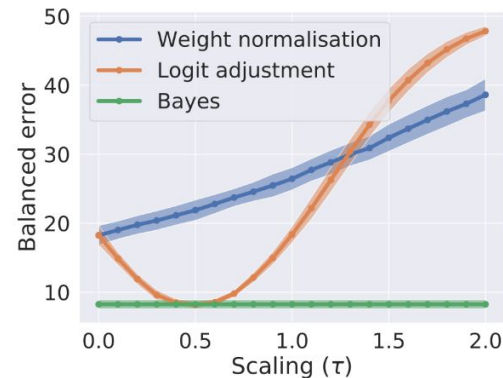
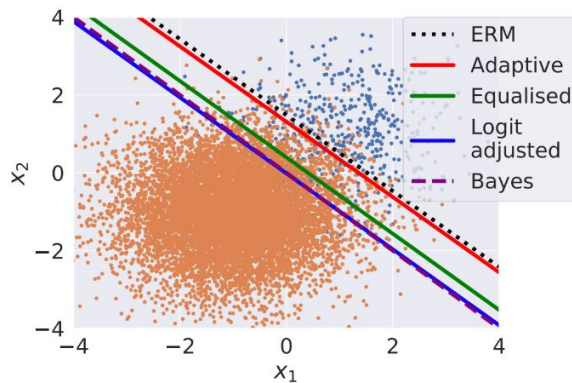
ERM will favour fewer mistakes on dominant class



Experiments: synthetic data



Converge to Bayes solution
(consistency)



Weight normalisation fails:
correct label has -ve score!

Experiments: real-world data

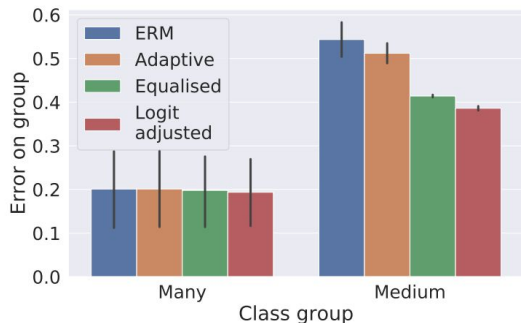
Method	CIFAR-10-LT	CIFAR-100-LT	ImageNet-LT	iNaturalist
ERM	27.16	61.64	53.11	38.66
Weight normalisation ($\tau = 1$) (Kang et al., 2020)	24.02	58.89	52.00	48.05
Weight normalisation ($\tau = \tau^*$) (Kang et al., 2020)	21.50	58.76	49.37	34.40*
Class-balanced (Cui et al., 2019)	25.43 [‡]	60.40 [‡]	53.21	35.84 [‡]
Adaptive (Cao et al., 2019)	26.65 [†]	60.40 [†]	52.15	33.31
Adaptive + DRW (Cao et al., 2019)	22.97 [†]	57.96 [†]	49.85	32.00 [†]
Equalised (Tan et al., 2020)	26.02	57.26	54.02	38.37
Logit adjustment post-hoc ($\tau = 1$)	22.60	58.24	49.66	33.98
Logit adjustment loss ($\tau = 1$)	22.33	56.11	48.89	33.64
Logit adjustment plus adaptive loss ($\tau = 1$)	22.42	55.92	51.25	31.56

Break-down of error rates

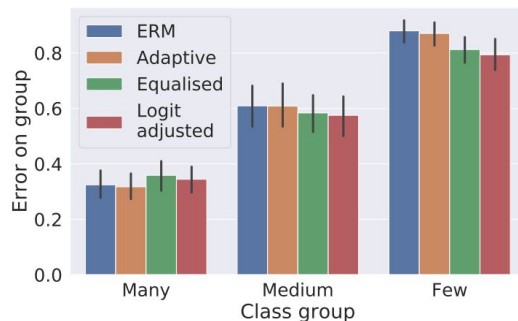
"Many": ≥ 100
examples

"Medium": [20, 100]
examples

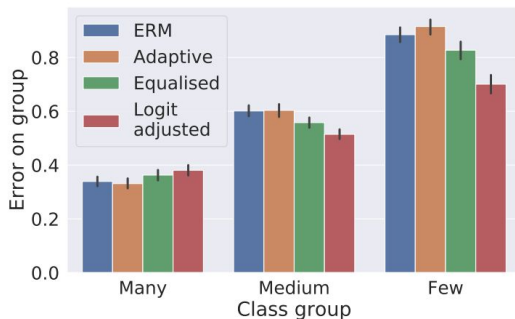
"Few": < 20
examples



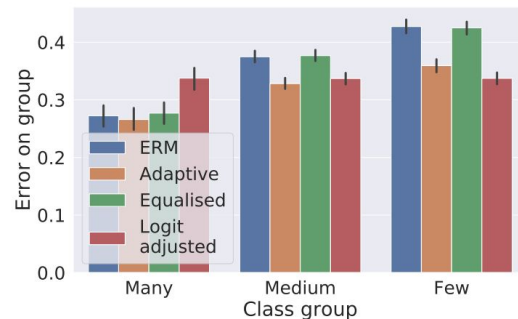
CIFAR10-LT



CIFAR100-LT



ImageNet-LT



iNaturalist

Summary

Summary of our work

- ✓ A **statistical perspective** of long-tail learning
- ✓ Unify and **generalise** existing works
- ✓ New **post-hoc** and **loss modification** approaches