

## Motivation

Consider a classification problem with  $K=3$  classes:  $\{Bird, Cat, Dog\}$ .



**One-vs-rest problem**

✓ Bird [ ] Cat [ ] Dog

**(I) Softmax classifier**

Use the prior #label=1.

Scale-up to a large number of classes.

Stable to train.

✗ Unable handle 0/N-label classification.

**Softmax activation:**

$$p_k = \frac{\exp(z_k)}{\sum_j \exp(z_j)}$$



**Multi-label problem**

[ ] Bird ✓ Cat ✓ Dog

**(II) Ensemble of binary classifiers**

Flexible to handle 0/N-label classification.

✗ Lack of correlation: independently trained.

✗ Do not scale-up to a large number of classes.

✗ Unstable to train.

**Sigmoid activation:**

$$p_k = \frac{\exp(z_k)}{1 + \exp(z_k)}$$



**Zero-label problem**

[ ] Bird [ ] Cat [ ] Dog

**Observation:** Binary classifier and softmax have similar form.

The difference is on the denominator (normalization factor).

## Method

**Goal:** Learn a normalization function  $C(X)$  that is shared across entire dataset  $X=\{x^{(i)}\}$  and classes.

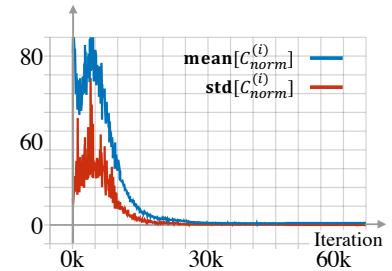
Probability for image  $x$  in being  $k$ -th class has the form:  $p(y = k|x) = q_k = \frac{\exp(z_k)}{C(X)}$ , that is constrained by:

- Prior knowledge  $\sum_{k=1}^K q_k = \text{\#label}$  for each image (See proof in paper.)
- Valid probability range  $q_k^{(i)} \in [0,1]$ .

**Training:** Maximize Likelihood Estimation, while

- 1) Computing  $C(X)$  over a mini-batch:  $\sum_i \sum_k \exp(z_k^{(i)}) / C = \sum_i \text{\#label}^{(i)}$
- 2) Minimizing  $-\max(\log q_k^{(i)}, 0)$  (Penalize any violation of  $q_k^{(i)} > 1$ ).

**Test:** Use the learned constant  $C$ :  $q_k = \min(\exp(z_k)/C, 1)$



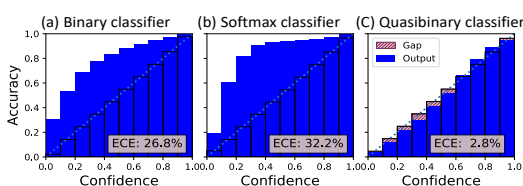
**$C(X)$  converges to a constant during training**

## Experiment

### Multi-label image classification

**Setup:** Follow Li *et al.* CVPR2017.

	MS-COCO		NUS-WIDE	
	$F_1 \uparrow$	ECE(%) $\downarrow$	$F_1 \uparrow$	ECE(%) $\downarrow$
Binary classifier [16], [18], [30]	51.2	26.8	40.7	23.6
Softmax [31]	<b>54.7</b>	32.2	43.2	25.8
Softmax w/ temperature [32]	<b>54.7</b>	31.4	43.2	24.6
Quasibinary classifier	<b>54.7</b>	<b>2.8</b>	<b>43.5</b>	<b>3.3</b>



1 Reliability diagrams for multi-label classification

**Conclusion:** Quasibinary classifier is both accurate and credible.

### One-vs.-rest image classification (single-label)

**Setup:** Resnet18 with 32x32 and 224x224 input.

	CIFAR10	CIFAR100	Tiny-ImageNet	ImageNet
Binary classifier	<b>4.8</b>	35.4	×	×
Quasibinary classifier (Ours)	4.9	<b>21.9</b>	<b>42.9</b>	25.4
Softmax classifier	5.2	22.2	43.3	<b>23.9</b>

**Conclusion:** Quasibinary classifier is better than binary classifier, and comparable with softmax classifier.

### Zero-label image classification

**Setup:** CIFAR60+40 dataset, with 40 classes of images from original CIFAR100 being treated as 0-label.

	IN Accuracy $\uparrow$	OUT MMC $\downarrow$	BOTH AU-ROC $\uparrow$
Binary classifiers [16], [18], [30]	77.8 %	14.7 %	0.901
Softmax + $\mathcal{L}_{\text{Uniform}}$ [6]	<b>80.7 %</b>	59.8 %	0.800
Softmax + $\mathcal{L}_{\text{MaxConf}}$ [11]	45.2 %	7.4 %	0.764
Quasibinary classifier	80.6 %	<b>6.9 %</b>	<b>0.913</b>

**Conclusion:** Quasibinary classifier achieves good performance on all measures.