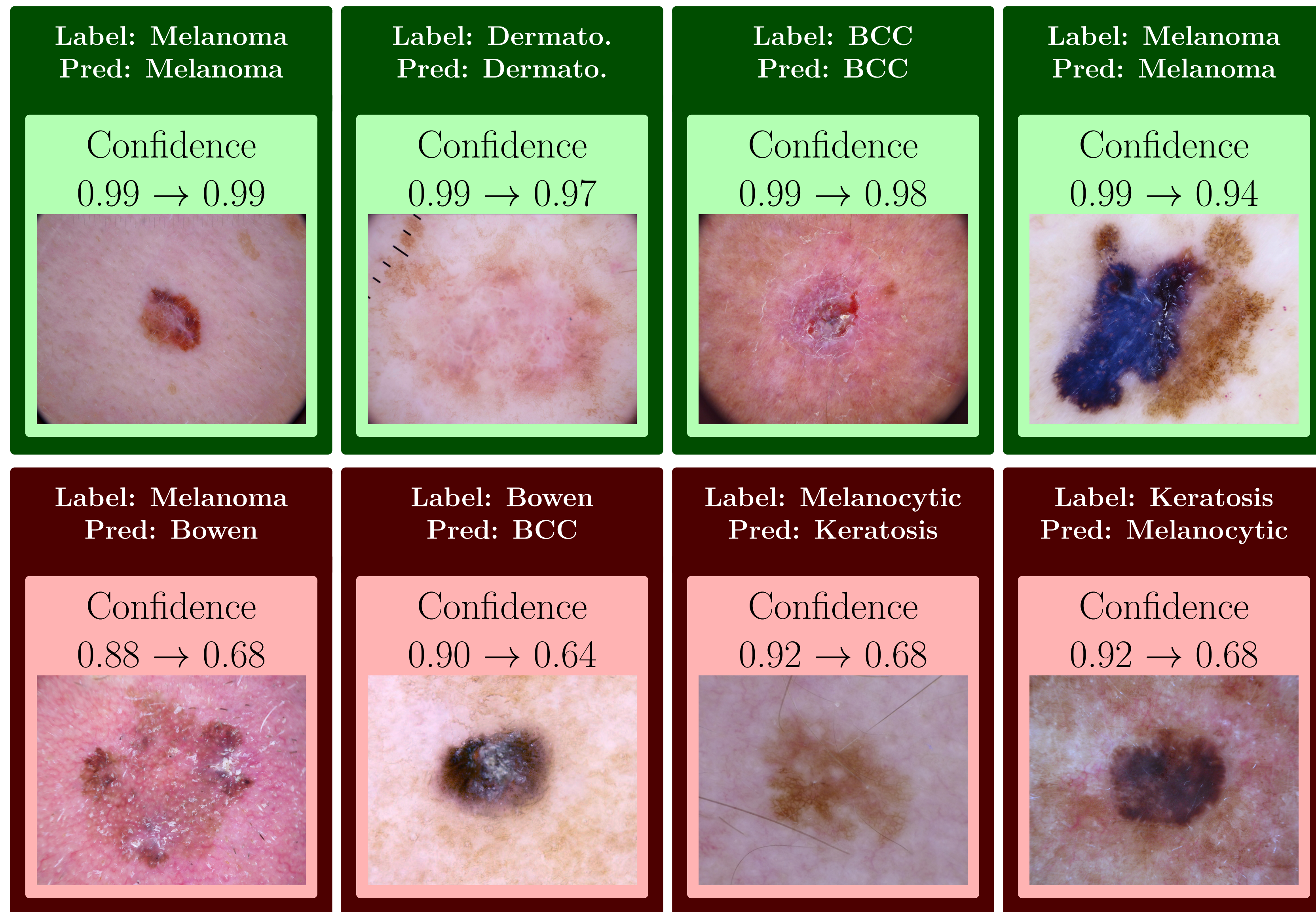


## Motivation

- Deep neural networks achieve impressive results but suffer from being overconfident
- Confidence is important to bring trust and reliability to decision-making applications



ISIC Dataset (Tschandl et al., 2018; Codella et al., 2017) ResNet152 Accuracy: 88.02%

## How to measure calibration?

### Negative Log-Likelihood (NLL)

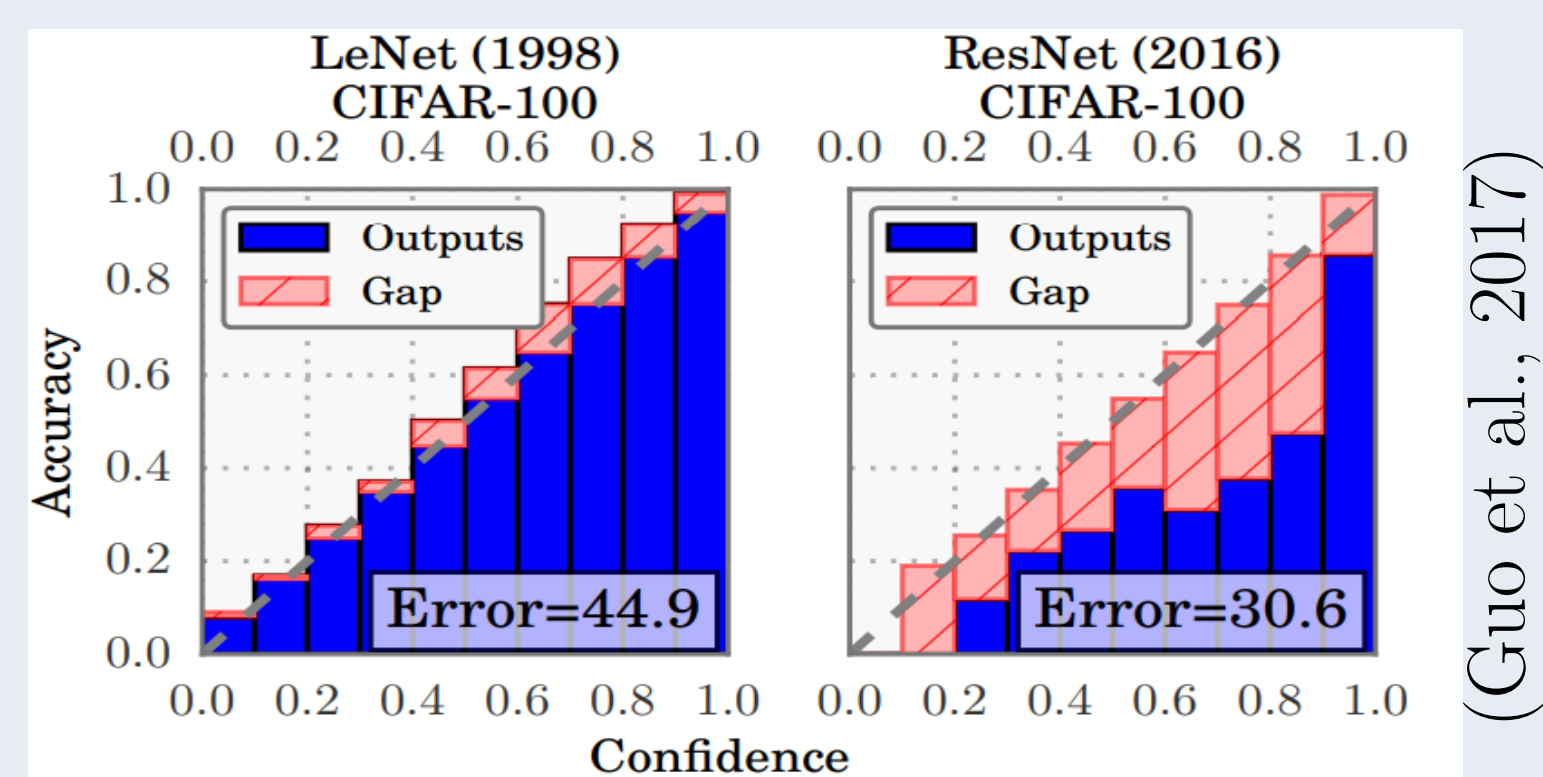
similarity between distributions  
softmax output  $S_y(\mathbf{x})$  and true conditional distribution  $Q(y|\mathbf{x})$

$$\text{NLL} = - \sum_{(\mathbf{x}_i, y_i)} \log \left( S_{y=y_i}(\mathbf{x}_i) \right), \quad (\mathbf{x}_i, y_i) \sim Q(\mathbf{x}, y)$$

### Expected Calibration Error (ECE)

x% of confidence  $\Rightarrow$  correct x% of the time

$$\text{ECE} = \sum_{l=1}^L \frac{|B_l|}{N} |\text{acc}(B_l) - \text{conf}(B_l)|$$



(Guo et al., 2017)

## Temperature Scaling (Guo et al., 2017)

- The goal is to rescale the logit layer with one single parameter  $T$  in order to minimize the calibration error of a pre-trained model.

$$S(h_i) = \frac{\exp^{h_i}}{\sum_j^K \exp^{h_j}} \Rightarrow S(h_i, T) = \frac{\exp^{\frac{h_i}{T}}}{\sum_j^K \exp^{\frac{h_j}{T}}}$$

$$T^* = \arg \min_T \left( \sum_{i=1}^N \overbrace{-\log \left( S_{y=y_i}(x_i, T) \right)}^{\text{NLL}} \right)$$

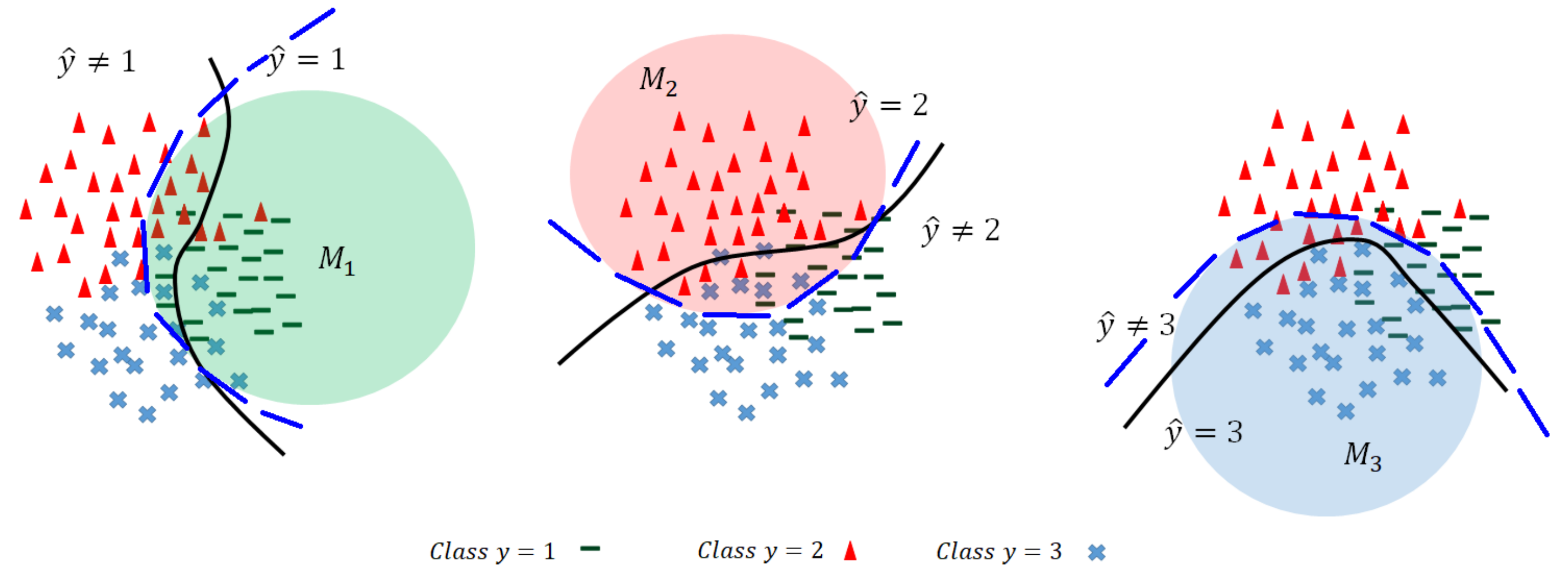
- Post-processing method
- Preserves accuracy
- Low time and memory complexity
- Highly sensitive to the size and labels of validation set
- Gathering enough correctly labelled validation data is costly

## Unsupervised Temperature Scaling

- How to make Temperature Scaling independent of labelled data?

### Key Insights:

- Estimate  $Q(y = k|\mathbf{x})$  for each class  $k$  instead of  $Q(y|\mathbf{x})$
- Use pre-trained model confidence to select the samples for each class distribution



### Approach:

- First, use confidence of pre-trained model  $S_{y=k}(\mathbf{x}_i)$  to select the samples of each class  $k$  distribution:

$$M_k = \{ \mathbf{x}_i \mid S_{y=k}(\mathbf{x}_i) \geq \theta_k, \quad \mathbf{x}_i \in \mathcal{V} \}$$

- How to select threshold  $\theta_k$ ?

$$\theta_k = \bar{S}_{y=k} + \sqrt{\frac{1}{|\mathcal{U}|} \sum_{\mathbf{x}_i \in \mathcal{U}} (\bar{S}_{y=k} - S_{y=k}(\mathbf{x}_i))^2}, \quad \text{where } \mathcal{U} = \{ \mathbf{x}_i \mid \hat{y}_i \neq k \}$$

- Finally, find  $T^*$  that minimize NLL for each class distribution:

$$\mathcal{L}_{UTS} = \sum_{k=1}^K \sum_{\mathbf{x}_i \in M_k} \overbrace{-\log \left( S_{y=k}(\mathbf{x}_i, T) \right)}^{\text{NLL}}$$

$$T^* = \arg \min_T (\mathcal{L}_{UTS})$$

## Analysis of samples' distribution selected in $M_k$

- Samples inside  $M_k$  are categorized in two groups:
  - The samples  $\hat{y} = k$  which are from  $Q(y = k|\mathbf{x})$
  - The samples  $\hat{y} \neq k$  which are located near to the decision boundary

From Bayes rule, we have:

$$Q(\mathbf{x}, y = k) = \frac{Q(y = k|\mathbf{x})}{Q(y \neq k|\mathbf{x})} Q(\mathbf{x}, y \neq k).$$

- Notice that  $Q(y = k|\mathbf{x}) \simeq Q(y \neq k|\mathbf{x})$  for samples located near to the decision boundary.
- Then, we can use samples generated from distributions  $Q(\mathbf{x}, y = k)$  or  $Q(\mathbf{x}, y \neq k)$  interchangeably.

## Experiments

### Benchmark Data

Model	Dataset	Accuracy	Uncalibrated		TS			UTS		
			NLL	ECE	NLL	ECE	$T$	NLL	ECE	$T$
DenseNet40	CIFAR10	92.61%	0.286	4.089	0.234	3.241	2.505	<b>0.221</b>	<b>0.773</b>	1.899
DenseNet40	CIFAR100	71.73%	1.088	8.456	<b>1.000</b>	<b>1.148</b>	1.450	1.001	1.945	1.493
DenseNet100	CIFAR10	95.06%	0.199	2.618	<b>0.156</b>	<b>0.594</b>	1.801	0.171	3.180	2.489
DenseNet100	CIFAR100	76.21%	1.119	11.969	0.886	4.742	2.178	<b>0.878</b>	<b>2.766</b>	1.694
DenseNet100	SVHN	95.72%	0.181	1.630	0.164	<b>0.615</b>	1.407	<b>0.162</b>	1.074	1.552
ResNet110	CIFAR10	93.71%	0.312	4.343	0.228	4.298	2.960	<b>0.207</b>	<b>1.465</b>	2.009
ResNet110	CIFAR100	70.31%	1.248	12.752	<b>1.051</b>	<b>1.804</b>	1.801	1.055	2.796	1.562
ResNet110	SVHN	96.06%	0.209	2.697	0.158	1.552	2.090	<b>0.152</b>	<b>0.550</b>	1.758
WideResNet32	CIFAR100	75.41%	1.166	13.406	0.909	<b>4.096</b>	2.243	<b>0.905</b>	4.872	1.651