

# On Focal Loss for Class-Posterior Probability Estimation: A Theoretical Perspective

Nontawat Charoenphakdee<sup>\*1,2</sup>, Jayakorn Vongkulbhisal<sup>\*3</sup>,  
Nuttapong Chairatanakul<sup>4,5</sup>, Masashi Sugiyama<sup>2,1</sup>

The University of Tokyo<sup>1</sup>, RIKEN AIP<sup>2</sup>, IBM Research<sup>3</sup>,  
Tokyo Institute of Technology<sup>4</sup>, RWBC-OIL (AIST)<sup>5</sup>

CVPR2021



**IBM Research**



# Multiclass classification

**Given:** input-output pairs:

$$\mathcal{Y} = \{1, 2, \dots, K\}$$

$e_y$ : One-hot vector

$$\{\mathbf{x}_i, y_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$$

We train a classifier  $g : \mathcal{X} \rightarrow \Delta^K$  by minimizing the empirical risk:

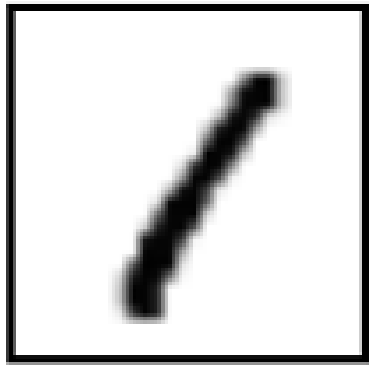
$$\hat{R}^\ell(g) = \frac{1}{n} \sum_{i=1}^n \ell(g(\mathbf{x}_i), e_{y_i}).$$

Loss function  $\ell$  highly influences the behavior of the trained classifier.

- A good classifier should predict the most probable class.
- But is this enough?

# Example

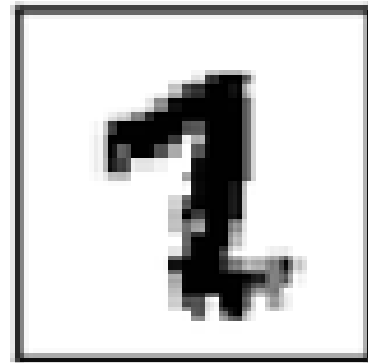
MNIST dataset



1

$p(y|\mathbf{x})$

0: 0.0  
**1: 1.0**  
 2: 0.0  
 ⋮  
 9: 0.0



1,2, or 7?

$p(y|\mathbf{x})$

0: 0.0  
**1: 0.3**  
**2: 0.5**  
 ⋮  
**7: 0.2**  
 ⋮  
 9: 0.0

Medical decision



Need surgery?

$p(y|\mathbf{x})$

$y = +1$   $\begin{pmatrix} .90 \\ .10 \end{pmatrix}$   
 $y = -1$

Class-posterior probability  $p(y|\mathbf{x})$  provides confidence score.

**Q: What loss can give  $p(y|\mathbf{x})$  ?**

# Cross-entropy loss

$$\ell_{\text{CE}}(\mathbf{v}, \mathbf{u}) = - \sum_{i=1}^K u_i \log(v_i)$$

$$\mathbf{u} \in \Delta^K, \mathbf{v} \in \Delta^K$$

\*  $\mathbf{u}$  is usually a one-hot label in practice

CE loss is **classification-calibrated**, i.e.,

CE risk minimizer gives the most probable class (Bayes-optimal):

$$\arg \max_y \mathbf{q}_{\text{CE}}^*(\mathbf{x}) = \arg \max_y p(y|\mathbf{x}).$$

CE loss is **strictly proper**, i.e.,

CE risk minimizer is a class-posterior probability estimator:

$$\mathbf{q}_{\text{CE}}^*(\mathbf{x}) = p(y|\mathbf{x}).$$

**Q: What about focal loss?**



Need surgery?

	$p(y \mathbf{x})$	$\mathbf{q}_{\text{CE}}^*(\mathbf{x})$
$y = +1$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$
$y = -1$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$

\*Strictly properness is sufficient to guarantee classification-calibration.

# Focal loss (Lin+, ICCV 2017)

$$\ell_{\text{FL}}^{\gamma}(\mathbf{v}, \mathbf{u}) = - \sum_{i=1}^K u_i (1 - v_i)^{\gamma} \log(v_i)$$

Originally proposed for dense object detection.

CE loss is a special case when  $\gamma = 0$ .

Focal loss has been used in many applications, e.g.,

- Electrocardiogram classification (Al Rahhal+, 2019)
- Brain tumor segmentation (Chang+, 2019)
- Femur fractures classification. (Lotfy+, 2019)

**Problem:** theoretical understanding of focal loss is limited.

**Q1: Is focal loss classification-calibrated?**

**Q2: Is focal loss strictly proper?**

# Main result

Focal loss is **classification-calibrated**:

$$\arg \max_y \mathbf{q}_{\text{FL},\gamma}^*(\mathbf{x}) = \arg \max_y p(y|\mathbf{x}).$$

However, it is **not strictly proper for  $\gamma > 0$** :

$$\mathbf{q}_{\text{FL},\gamma}^*(\mathbf{x}) \neq p(y|\mathbf{x}).$$



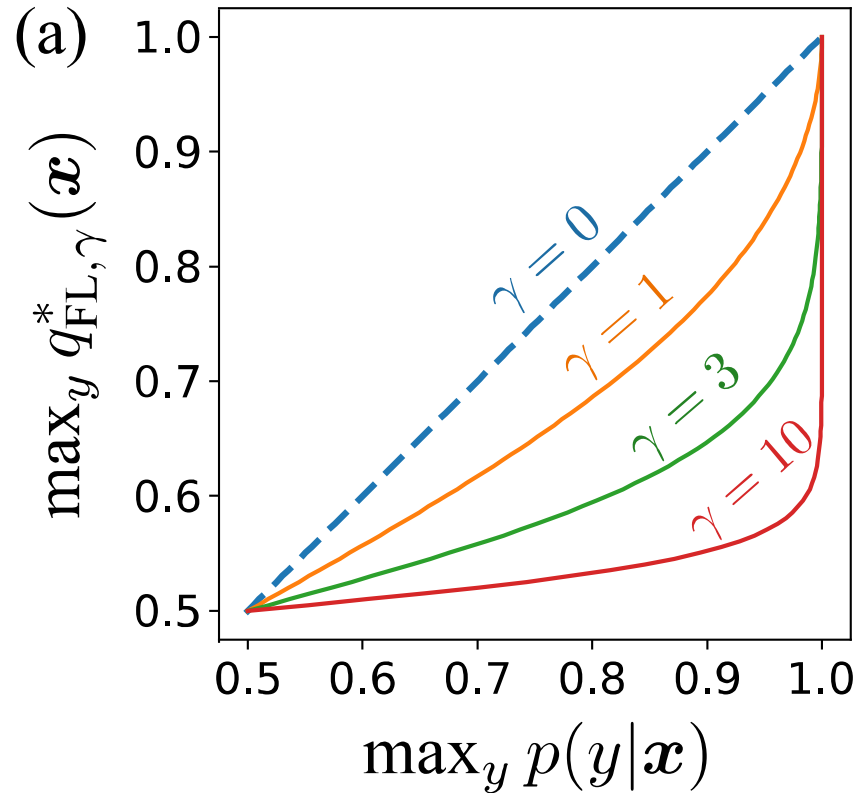
Need surgery?

	$p(y \mathbf{x})$	$\mathbf{q}_{\text{CE}}^*(\mathbf{x})$	$\mathbf{q}_{\text{FL},1}^*(\mathbf{x})$	$\mathbf{q}_{\text{FL},3}^*(\mathbf{x})$	$\mathbf{q}_{\text{FL},5}^*(\mathbf{x})$
$y = +1$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .78 \\ .22 \end{pmatrix}$	$\begin{pmatrix} .65 \\ .35 \end{pmatrix}$	$\begin{pmatrix} .60 \\ .40 \end{pmatrix}$
$y = -1$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .22 \\ .78 \end{pmatrix}$	$\begin{pmatrix} .35 \\ .65 \end{pmatrix}$	$\begin{pmatrix} .40 \\ .60 \end{pmatrix}$

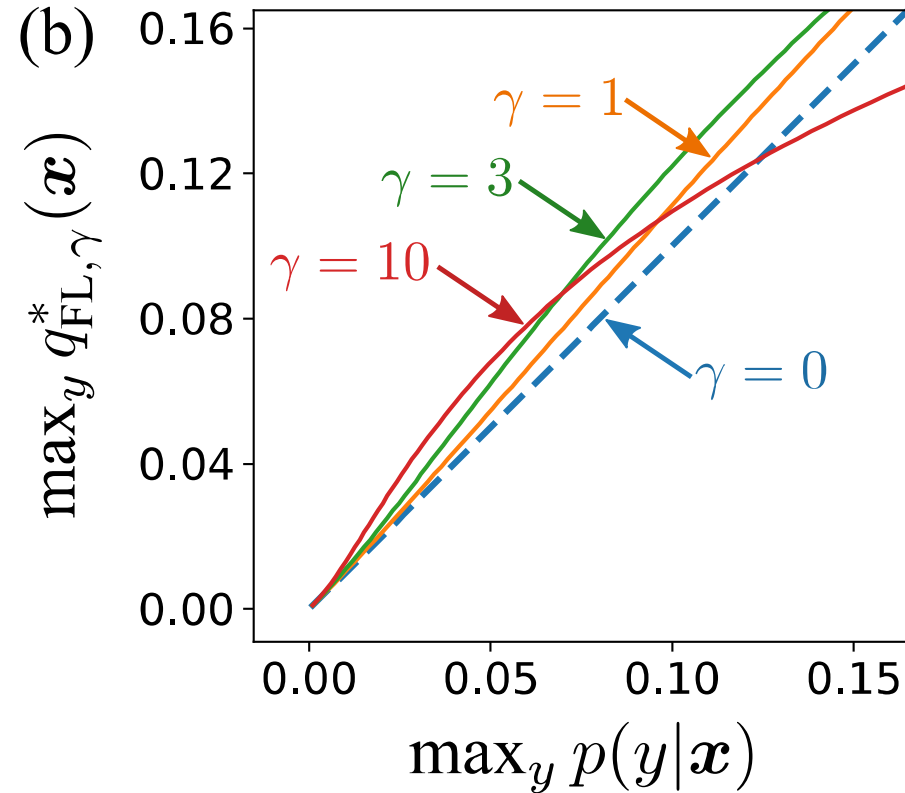
We can predict the most probable class, but **confidence score is unreliable!**

# Focal risk minimizer can be both under/overconfident

Underconfident (K=2)



Overconfident (K=1000)



**Q: How to solve this problem?**

*(Please see our paper for more detail.)*

# Solution: Recover $p(y|\mathbf{x})$ from $q_{\text{FL},\gamma}^*(\mathbf{x})$ via $\Psi^\gamma$



Need surgery?

	$p(y \mathbf{x})$	$q_{\text{CE}}^*(\mathbf{x})$	$q_{\text{FL},1}^*(\mathbf{x})$	$q_{\text{FL},3}^*(\mathbf{x})$	$q_{\text{FL},5}^*(\mathbf{x})$
$y = +1$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$	$\begin{pmatrix} .78 \\ .22 \end{pmatrix}$	$\begin{pmatrix} .65 \\ .35 \end{pmatrix}$	$\begin{pmatrix} .60 \\ .40 \end{pmatrix}$
$y = -1$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .10 \\ .90 \end{pmatrix}$	$\begin{pmatrix} .22 \\ .78 \end{pmatrix}$	$\begin{pmatrix} .35 \\ .65 \end{pmatrix}$	$\begin{pmatrix} .40 \\ .60 \end{pmatrix}$

$\Psi^1 \downarrow$

$$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$$

$\Psi^3 \downarrow$

$$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$$

$\Psi^5 \downarrow$

$$\begin{pmatrix} .90 \\ .10 \end{pmatrix}$$

Define  $\Psi^\gamma(\mathbf{v}) = [\Psi_1^\gamma(\mathbf{v}), \dots, \Psi_K^\gamma(\mathbf{v})]^\top$ ,

where  $\Psi_i^\gamma(\mathbf{v}) = \frac{h^\gamma(v_i)}{\sum_{l=1}^K h^\gamma(v_l)}$ ,

and  $h^\gamma(v_i) = \frac{v_i}{(1-v_i)^\gamma - \gamma(1-v_i)^{\gamma-1} v_i \log v_i}$ .

- Closed-form
- No hyperparameter
- Theoretically justified
- Preserves accuracy
- No additional training required

(Please see our paper for experiments on benchmark datasets.)



# Conclusions

**Theoretical analysis of focal loss with practical use.**

Q1: Is focal loss classification-calibrated?

**Yes!**

Q2: Is focal loss strictly proper?

**No!** Directly using model's output gives **unreliable confidence**.

Q3: Following Q2, can we do anything about it?

**Yes!** We discovered a closed-form transformation  $\Psi^\gamma$  that can recover  $p(y|\mathbf{x})$  **with theoretical guarantee!**