



# Generalizing Consistent Multi-Class Classification with Rejection to be Compatible with Arbitrary Losses

**Yuzhou Cao<sup>1\*</sup> Tianchi Cai<sup>2</sup> Lei Feng<sup>3</sup> Lihong Gu<sup>2</sup> Jinjie Gu<sup>2</sup>  
Bo An<sup>1</sup> Gang Niu<sup>4</sup> Masashi Sugiyama<sup>4,5</sup>**

<sup>1</sup>Nanyang Technological University    <sup>2</sup>Ant Group

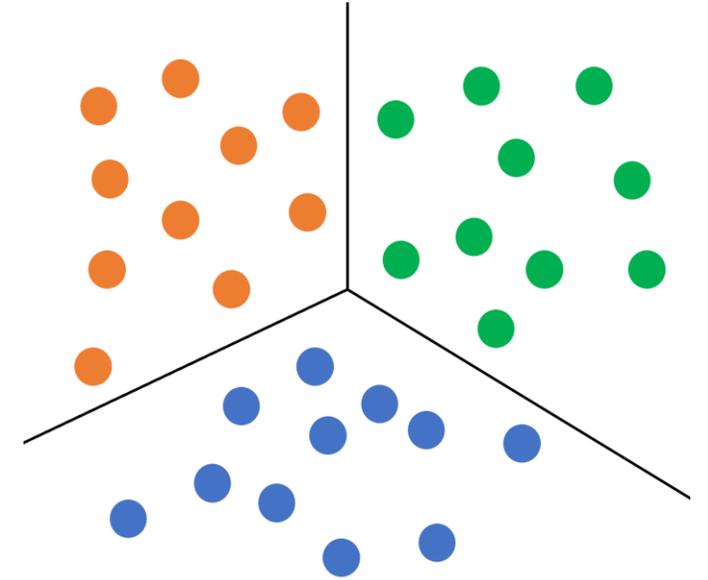
<sup>3</sup>Chongqing University    <sup>4</sup>RIKEN

<sup>5</sup>The University of Tokyo

# Introduction

## ➤ Multi-Class Classification

Learn a classifier to decide the exact class label of each data point.



# Introduction

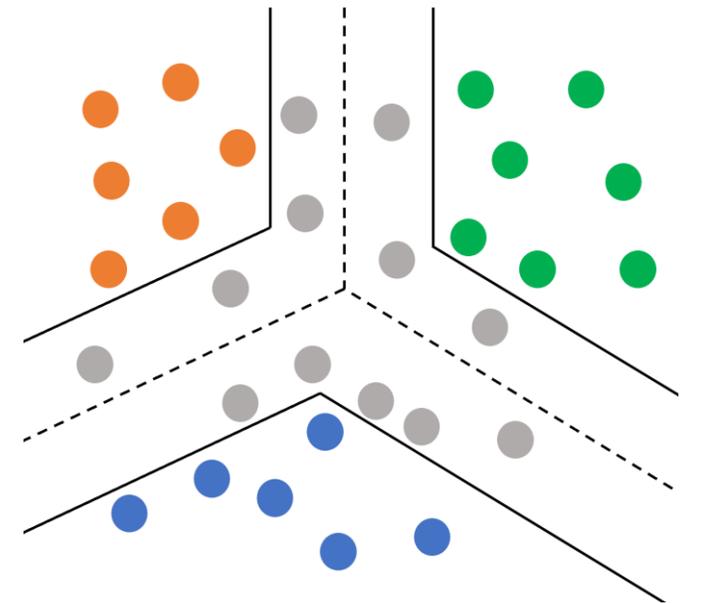
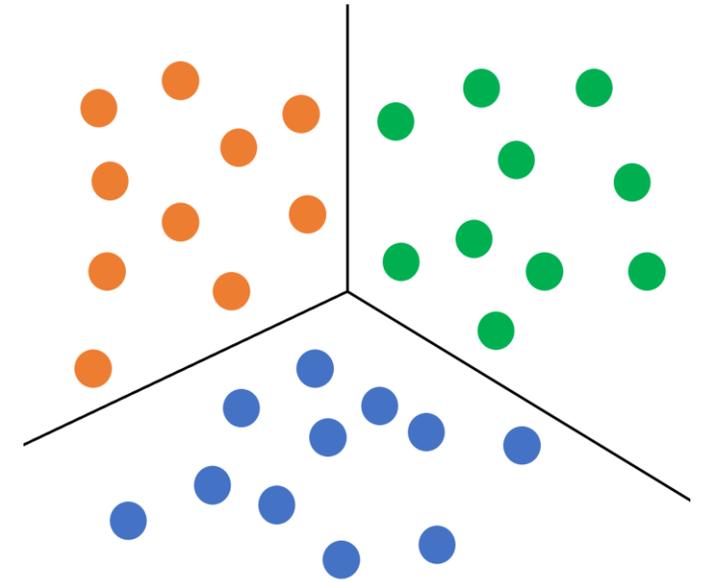
## ➤ Multi-Class Classification

Learn a classifier to decide the exact class label of each data point.



## ➤ Classification with Rejection

Learn a classifier and rejector that classifies **safe** samples and rejects **risky** samples.



# Introduction

## ➤ Multi-Class Classification

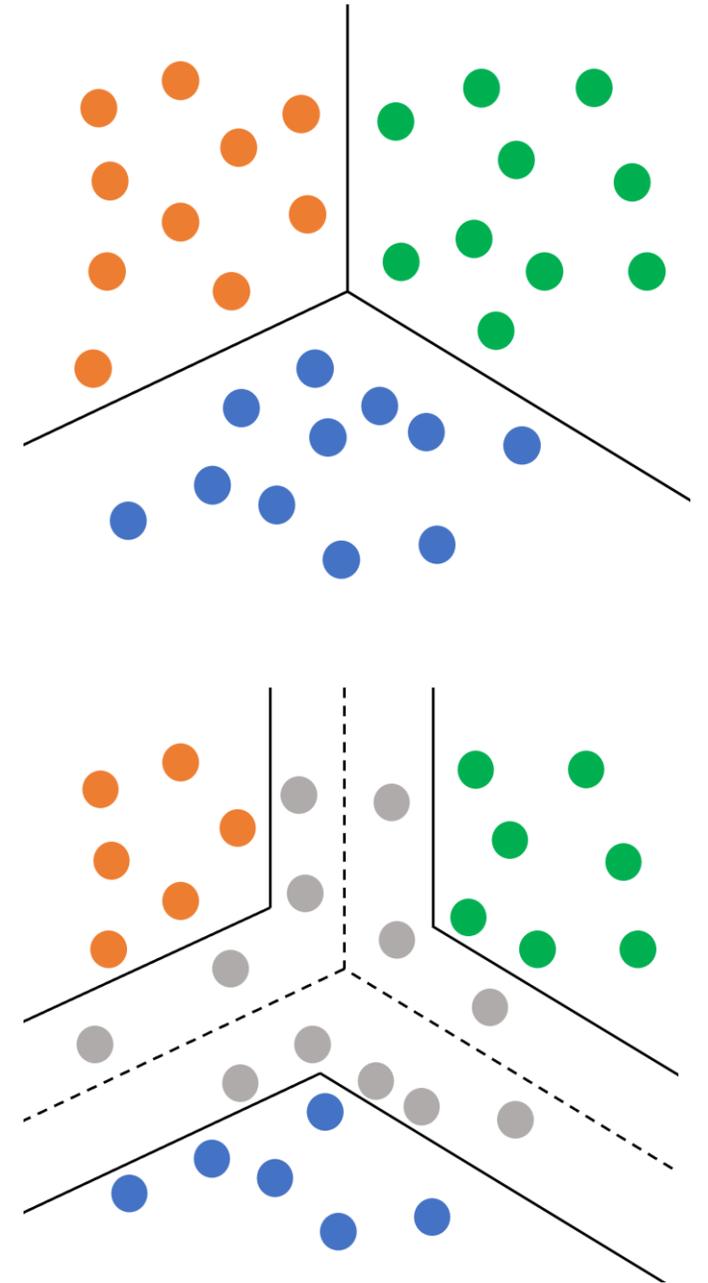
Learn a classifier to decide the exact class label of each data point.



## ➤ Classification with Rejection

Learn a classifier and rejector that classifies **safe** samples and rejects **risky** samples.

**What does 'risky' means?**

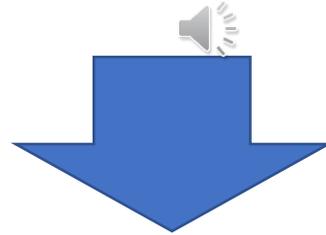


# Hard Samples are Risky to be Classified

## ➤ Examples from MNIST

$$p(\text{7} = \text{digit '7'}) = 50\% \quad p(\text{7} = \text{digit '3'}) = 50\%$$

$$p(\text{6} = \text{digit '5'}) = 50\% \quad p(\text{6} = \text{digit '6'}) = 50\%$$



Minimal misclassification rate: 50%!

# Optimal Rejection Rule

## ➤ Chow's Rule [1]

$$f^{chow}(\mathbf{x}) = \begin{cases} \operatorname{argmax}_y p(y|\mathbf{x}), & \text{if } \max_y p(y|\mathbf{x}) > 1 - c, \\ \text{reject}, & \text{else.} \end{cases}$$

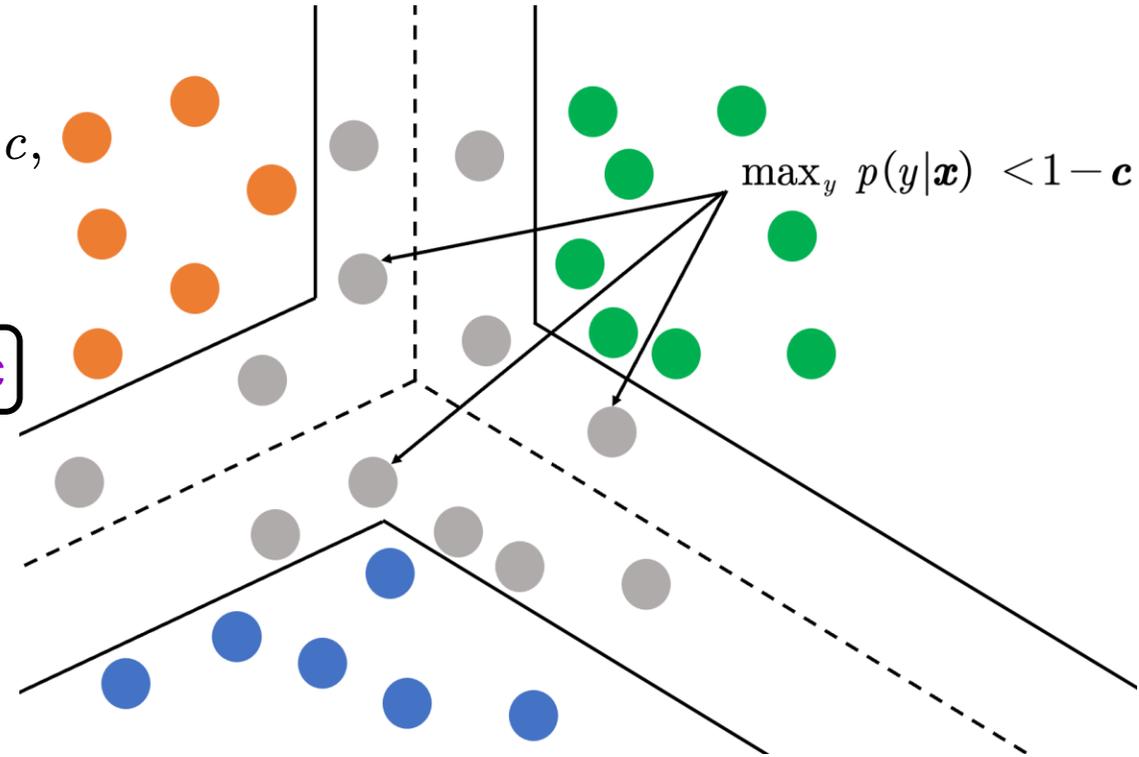


# Optimal Rejection Rule

## ➤ Chow's Rule [1]

$$f^{chow}(\mathbf{x}) = \begin{cases} \operatorname{argmax}_y p(y|\mathbf{x}), & \text{if } \max_y p(y|\mathbf{x}) > 1 - c, \\ \text{reject}, & \text{else.} \end{cases}$$

Rejecting samples with misclassification rate  $> c$



# Optimal Rejection Rule

## ➤ Chow's Rule [1]

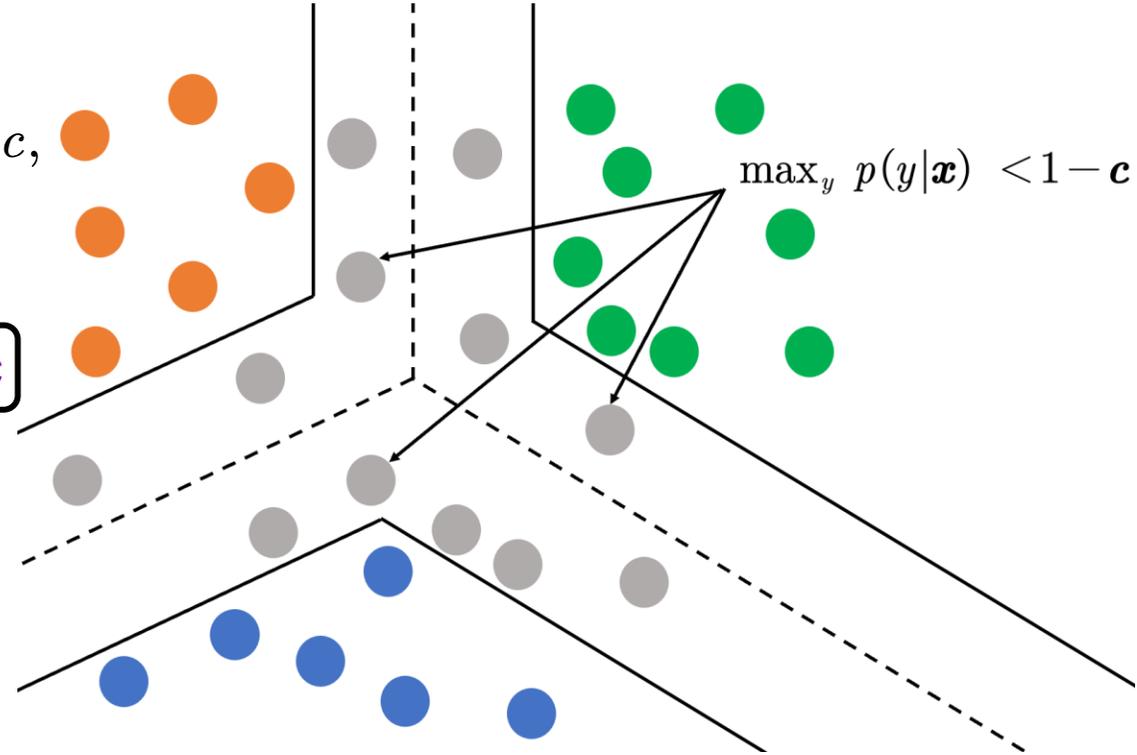
$$f^{chow}(\mathbf{x}) = \begin{cases} \operatorname{argmax}_y p(y|\mathbf{x}), & \text{if } \max_y p(y|\mathbf{x}) > 1 - c, \\ \text{reject}, & \text{else.} \end{cases}$$

Rejecting samples with misclassification rate  $> c$

## ➤ Risk Minimization Framework:

$$R_{01c}(f) = \mathbb{E}_{p(\mathbf{x}, y)} [\ell_{01c}(f(\mathbf{x}), y)]$$

$$\text{where } \ell_{01c}(f(\mathbf{x}), y) = \begin{cases} c, & f(\mathbf{x}) = \text{reject}, \\ \mathbb{I}(f(\mathbf{x}) \neq y), & \text{else.} \end{cases}$$



# Optimal Rejection Rule

## ➤ Chow's Rule [1]

$$f^{chow}(\mathbf{x}) = \begin{cases} \operatorname{argmax}_y p(y|\mathbf{x}), & \text{if } \max_y p(y|\mathbf{x}) > 1 - c, \\ \text{reject}, & \text{else.} \end{cases}$$

Rejecting samples with misclassification rate  $> c$

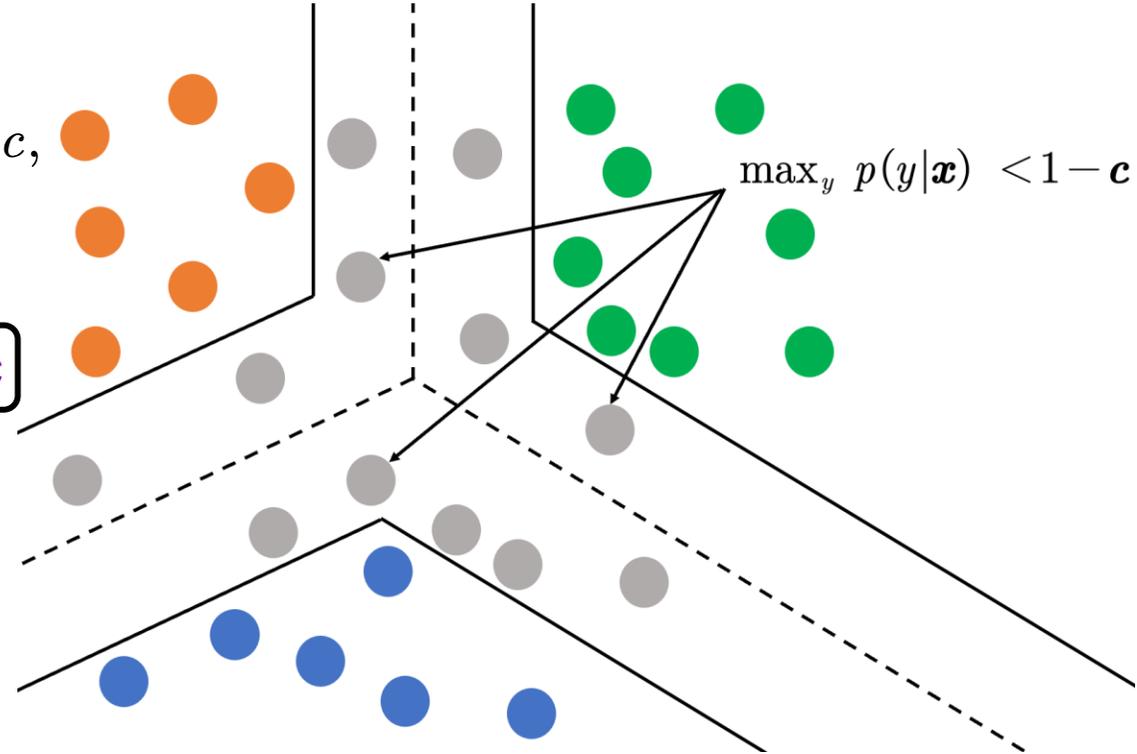
## ➤ Risk Minimization Framework:

$$R_{01c}(f) = \mathbb{E}_{p(\mathbf{x}, y)} [\ell_{01c}(f(\mathbf{x}), y)]$$

$$\text{where } \ell_{01c}(f(\mathbf{x}), y) = \begin{cases} c, & f(\mathbf{x}) = \text{reject}, \\ \mathbb{I}(f(\mathbf{x}) \neq y), & \text{else.} \end{cases}$$



$$\operatorname{argmin}_f R_{01c}(f) = f^{chow}$$



# Optimal Rejection Rule

## ➤ Chow's Rule [1]

$$f^{chow}(\mathbf{x}) = \begin{cases} \operatorname{argmax}_y p(y|\mathbf{x}), & \text{if } \max_y p(y|\mathbf{x}) > 1 - c, \\ \text{reject}, & \text{else.} \end{cases}$$

Rejecting samples with misclassification rate  $> c$

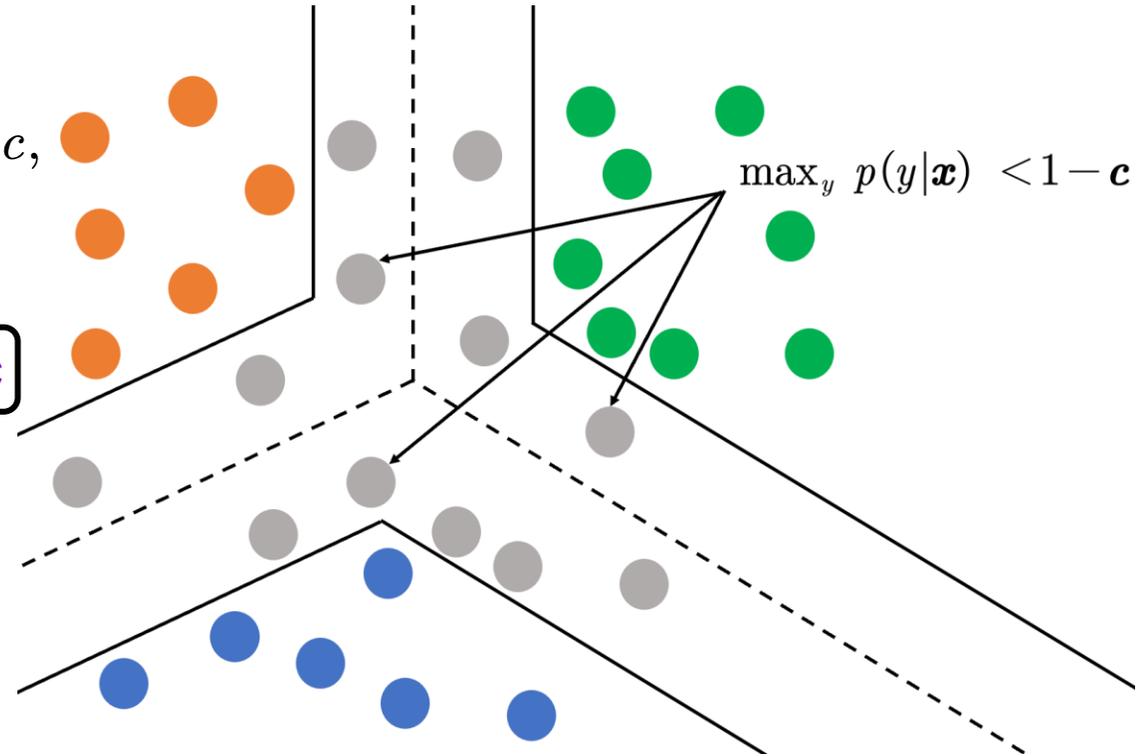
## ➤ Risk Minimization Framework:

$$R_{01c}(f) = \mathbb{E}_{p(\mathbf{x}, y)} [\ell_{01c}(f(\mathbf{x}), y)]$$

$$\text{where } \ell_{01c}(f(\mathbf{x}), y) = \begin{cases} c, & f(\mathbf{x}) = \text{reject}, \\ \mathbb{I}(f(\mathbf{x}) \neq y), & \text{else.} \end{cases}$$

$$\operatorname{argmin}_f R_{01c}(f) = f^{chow}$$

Intractable! (NP-hard)



# Calibrated Surrogate Losses

- **Class-Posterior Probability Estimation Based [2]:**

$\operatorname{argmax}_y \hat{p}(y|\mathbf{x}) \rightarrow 1$  **Severe overconfidence!**



# Calibrated Surrogate Losses

- **Class-Posterior Probability Estimation Based [2]:**

$\operatorname{argmax}_y \hat{p}(y|\mathbf{x}) \rightarrow 1$  **Severe overconfidence!**

- **Cost-Sensitive Learning Based [3], Learning to Defer [4]:**

**Restriction on loss functions (OvA Loss, CE Loss)**

# Problem Reduction: K+1 Class Classification

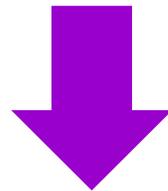
➤ A New Data Generation Distribution:

$$\tilde{p}(\mathbf{x}, \tilde{y}) = \begin{cases} \frac{p(\mathbf{x}, y)}{2-c}, & \tilde{y} \in \{1, 2, \dots, K\}, \\ \frac{(1-c)p(\mathbf{x})}{2-c}, & \tilde{y} = \text{rejected}. \end{cases}$$

# Problem Reduction: K+1 Class Classification

➤ A New Data Generation Distribution:

$$\tilde{p}(\mathbf{x}, \tilde{y}) = \begin{cases} \frac{p(\mathbf{x}, y)}{2-c}, & \tilde{y} \in \{1, 2, \dots, K\}, \\ \frac{(1-c)p(\mathbf{x})}{2-c}, & \tilde{y} = \text{rejected}. \end{cases}$$



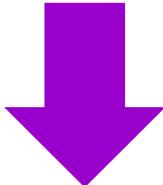
$$\tilde{R}_{01}(f) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})} [\mathbb{I}(f(\mathbf{x}) \neq \tilde{y})]$$

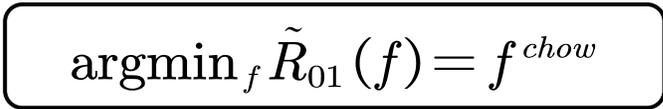
$$\operatorname{argmin}_f \tilde{R}_{01}(f) = f^{chow}$$

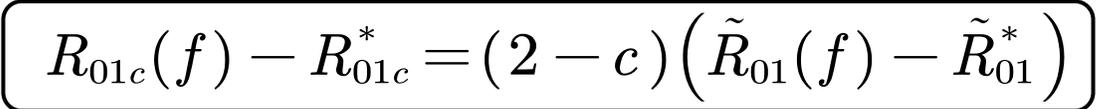
# Problem Reduction: K+1 Class Classification

➤ A New Data Generation Distribution:

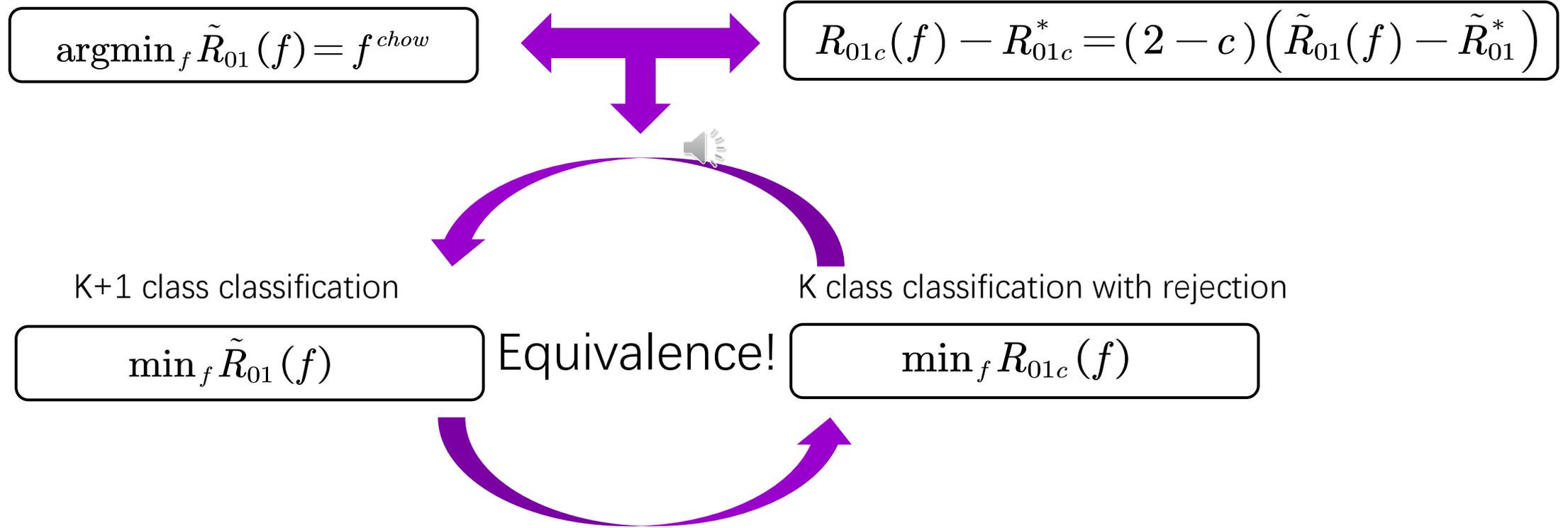
$$\tilde{p}(\mathbf{x}, \tilde{y}) = \begin{cases} \frac{p(\mathbf{x}, y)}{2-c}, & \tilde{y} \in \{1, 2, \dots, K\}, \\ \frac{(1-c)p(\mathbf{x})}{2-c}, & \tilde{y} = \text{rejected}. \end{cases}$$


$$\tilde{R}_{01}(f) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})} [\mathbb{I}(f(\mathbf{x}) \neq \tilde{y})]$$


$$\operatorname{argmin}_f \tilde{R}_{01}(f) = f^{chow}$$


$$R_{01c}(f) - R_{01c}^* = (2-c) (\tilde{R}_{01}(f) - \tilde{R}_{01}^*)$$

# Problem Reduction: K+1 Class Classification



# Problem Reduction: K+1 Class Classification

- **How to  $\min_f \tilde{R}_{01}(f)$ ? Not trivial!**
  - ①. No data from the class 'rejected'!
  - ②. Optimization of 0-1 loss is NP-hard!



# Problem Reduction: K+1 Class Classification

➤ How to  $\min_f \tilde{R}_{01}(f)$ ? **Not trivial!**

①. No data from the class 'rejected'!

②. Optimization of 0-1 loss is NP-hard!

➤  $\min_g R_{L_c^\Phi}(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)} [L_\Phi(\mathbf{g}(\mathbf{x}), y)]$  **instead!**

where  $L_c^\Phi(\mathbf{u}, y) = \Phi(\mathbf{u}, y) + (1 - c)\Phi(\mathbf{u}, K + 1)$   $f(\mathbf{x}) = \begin{cases} \text{reject, } \operatorname{argmax}_y \mathbf{u}_y = K + 1 \\ \operatorname{argmax}_y \mathbf{u}_y, \text{ else.} \end{cases}$

➤ **Motivation:**  $R_{L_c^\Phi}(\mathbf{g}) = (2 - c) \tilde{R}_\Phi(\mathbf{g})$

where  $\tilde{R}_\Phi(\mathbf{g}) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})} [\Phi(\mathbf{g}(\mathbf{x}), \tilde{y})]$

# Problem Reduction: K+1 Class Classification

➤ How to  $\min_f \tilde{R}_{01}(f)$ ? Not trivial!

①. No data from the class 'rejected'!

②. Optimization of 0-1 loss is NP-hard!

➤  $\min_g R_{L_c^\Phi}(\mathbf{g}) = \mathbb{E}_{p(\mathbf{x}, y)} [L_\Phi(\mathbf{g}(\mathbf{x}), y)]$  instead!

where  $L_c^\Phi(\mathbf{u}, y) = \Phi(\mathbf{u}, y) + (1 - c)\Phi(\mathbf{u}, K + 1)$   $f(\mathbf{x}) = \begin{cases} \text{reject, } \operatorname{argmax}_y \mathbf{u}_y = K + 1 \\ \operatorname{argmax}_y \mathbf{u}_y, \text{ else.} \end{cases}$

➤ **Motivation:**  $R_{L_c^\Phi}(\mathbf{g}) = (2 - c) \tilde{R}_\Phi(\mathbf{g})$

where  $\tilde{R}_\Phi(\mathbf{g}) = \mathbb{E}_{\tilde{p}(\mathbf{x}, \tilde{y})} [\Phi(\mathbf{g}(\mathbf{x}), \tilde{y})]$

**Main Result:**

*Any* classification-calibrated surrogate  $\Phi$  can make  $L_c^\Phi$  calibrated w.r.t.  $\ell_{01c}$ .

$$R_{L_c^\Phi}(\mathbf{g}_i) \rightarrow R_{L_c^\Phi}^* \Rightarrow R_{01c}(\operatorname{argmax}_i \mathbf{g}_i) \rightarrow R_{01c}^*$$

# More Discoveries

- Regret Transfer Bound,
- Estimation Error Bound,
- Analysis of GCE Loss,
- Instance-Dependent Cost
- Experimental Results...



# References

- [1]. C. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*, 16(1):41–46, 1970.
- [2]. Chenri Ni, Nontawat Charoenphakdee, Junya Honda, and Masashi Sugiyama. On the calibration of multiclass classification with rejection. In *NeurIPS*, 2019. 
- [3]. Nontawat Charoenphakdee, Zhenghang Cui, Yivan Zhang, and Masashi Sugiyama. Classification with rejection based on cost-sensitive classification. In *ICML*, 2021.
- [4]. Hussein Mozannar and David A. Sontag. Consistent estimators for learning to defer to an expert. In *ICML*, 2020.