

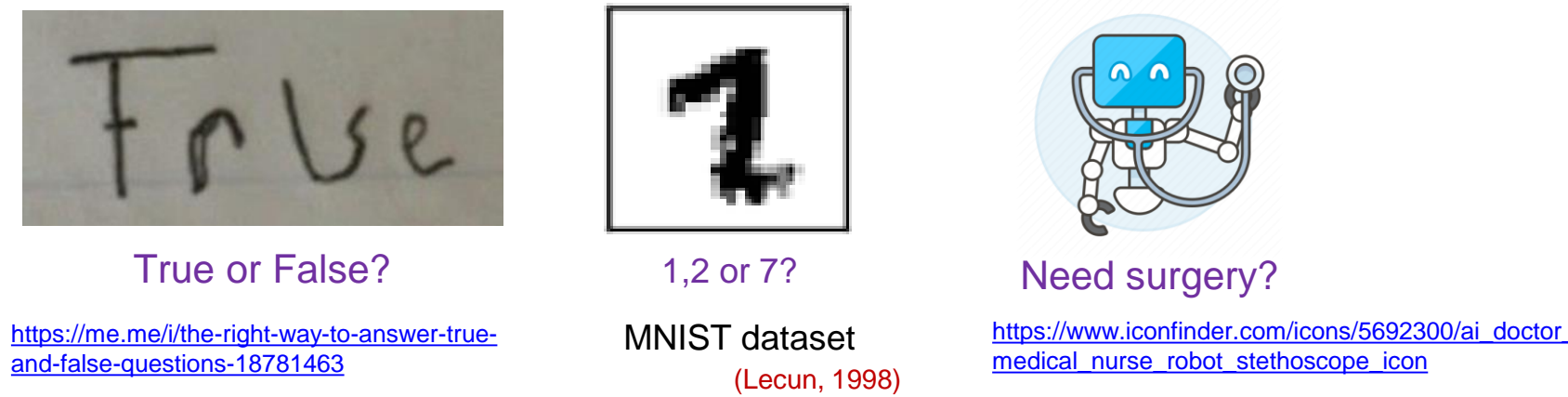
Classification with Rejection Based on Cost-sensitive Classification

Nontawat Charoenphakdee^{1,2}, Zhenghang Cui^{1,2}, Yivan Zhang^{1,2}, Masashi Sugiyama^{2,1}

¹The University of Tokyo, ²RIKEN AIP

International conference on machine learning (ICML) 2021

Summary



Saying “I don’t know” can **prevent misclassification**.
How to learn a classifier to say “I don’t know” reasonably?

The well-known confidence-based approach typically **requires estimating** $p(y|\mathbf{x})$.
Theoretical framework typically **requires a loss to be convex**. (Ni et al., 2019)
Existing approaches have **less loss choice than that of ordinary classification**.

Contributions:

We propose a cost-sensitive approach to classification with rejection.

1. It can **avoid estimating** $p(y|\mathbf{x})$.
2. It is applicable to **both binary and multiclass cases**.
3. It is **theoretically justifiable for any classification-calibrated loss***

*Classification calibration is known to be a minimum loss requirement for ordinary classification.

Problem formulation

Given: rejection cost $c \in (0, 0.5)$, training input-output pairs:

$$\{\mathbf{x}_i, y_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y) \quad \mathcal{Y} = \{1, 2, \dots, K\}: \text{Label space}$$

$$f: \mathcal{X} \rightarrow \mathcal{Y} \cup \{\textcircled{R}\}: \text{Classification rule}$$

Goal: minimize 0-1-c risk:

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

$$\ell_{01c}(f(\mathbf{x}), y) = \begin{cases} c & \text{if } f(\mathbf{x}) = \textcircled{R} \\ \ell_{01}(f(\mathbf{x}), y) & \text{otherwise} \end{cases}$$

0-1 loss

Rejection cost c is less than **misclassification cost**.

Directly minimizing the empirical 0-1-c risk is computationally infeasible. (Bartlett and Wegkamp, 2008)

Bayes-optimal solution: Chow’s rule

Knowing $p(y|\mathbf{x})$ is sufficient to obtain optimal solution. (Chow 1970)

$$f^*(\mathbf{x}) = \begin{cases} \textcircled{R} & \max_y p(y|\mathbf{x}) \leq 1 - c, \\ \arg \max_y p(y|\mathbf{x}) & \text{otherwise.} \end{cases}$$

Straightforward solution: estimating $p(y|\mathbf{x})$ (confidence-based approach).

- **More restrictive loss requirement** than classification calibration. (Reid and Williamson, 2010)

Well-known loss such as hinge, ramp, and sigmoid losses are classification-calibrated but not capable of estimating $p(y|\mathbf{x})$.

Q: Can we have a framework that can use any classification-calibrated loss?

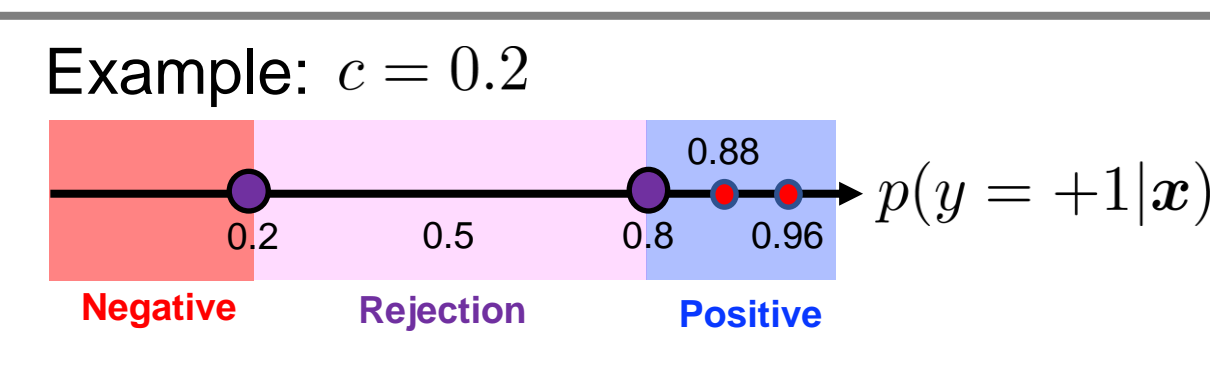
Proposal: Cost-sensitive approach

Binary case

Chow’s rule for the binary case: $f^*(\mathbf{x}) = \begin{cases} +1 & p(y=1|\mathbf{x}) > 1 - c, \\ \textcircled{R} & c \leq p(y=1|\mathbf{x}) \leq 1 - c, \\ -1 & p(y=1|\mathbf{x}) < c, \end{cases}$

To mimic Chow’s rule, we only need to know:

1. $p(y=1|\mathbf{x}) > 1 - c$
2. $p(y=1|\mathbf{x}) < c$



Binary cost-sensitive classification:

Binary classification where false negative penalty \neq false positive penalty.
Let false positive penalty be $\alpha \in (0, 1)$ and false negative penalty be $1 - \alpha$:

- Solving cost-sensitive classification can validate if $p(y=1|\mathbf{x}) > \alpha$.
- Loss requirement: **classification calibration** (Scott, 2012)

Proposal: cost-sensitive approach to binary classification with rejection.

- Learn two cost-sensitive classifiers for $\alpha = c$ and $\alpha = 1 - c$.
- Predict if both classifiers predict the same class and reject otherwise.

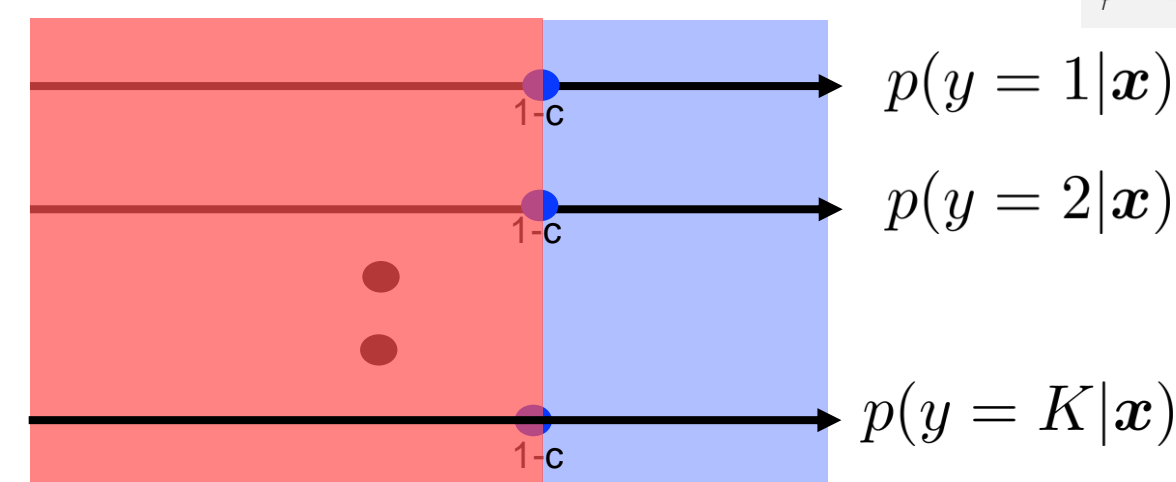
Multiclass extension

Learn K **one-vs-rest cost-sensitive classifiers** with $\alpha = 1 - c$.

Can be learned at once by learning $g: \mathcal{X} \rightarrow \mathbb{R}^K$.

$$\mathcal{L}_{CS}^{c, \phi}(g; \mathbf{x}, y) = c\phi(g_y(\mathbf{x})) + (1 - c) \sum_{y' \neq y} \phi(-g_{y'}(\mathbf{x})).$$

$\phi: \mathbb{R} \rightarrow \mathbb{R}$: Classification-calibrated margin loss



- Predict if:**
Only one classifier returns positive
- Reject if:**
1. All classifiers predict negative, or
2. More than one classifier predicts positive

Classification rule:

$$f(\mathbf{x}; g) = \begin{cases} \textcircled{R} & \max_y g_y(\mathbf{x}) \leq 0, \\ \textcircled{R} & \exists y, y' \text{ s.t. } y \neq y' \\ & g_y(\mathbf{x}), g_{y'}(\mathbf{x}) > 0, \\ \arg \max_y g_y(\mathbf{x}) & \text{otherwise.} \end{cases}$$

Excess risk bound

$$\text{Main result: } \underbrace{R^{\ell_{01c}}(f) - R^{\ell_{01c},*}}_{\text{Excess 0-1-c risk}} \leq \underbrace{R^{\mathcal{L}_{CS}^{c, \ell_{01}}}(g) - R^{\mathcal{L}_{CS}^{c, \ell_{01},*}}}_{\text{Excess cost-sensitive 0-1 risk}}$$

Excess 0-1-c risk is bounded by excess cost-sensitive 0-1 risk!

Excess risk bound of cost-sensitive 0-1 risk is well studied. (Scott 2012, Steinwart, 2007)

$$R^{\mathcal{L}_{CS}^{c, \ell_{01}}}(g) - R^{\mathcal{L}_{CS}^{c, \ell_{01},*}} \leq \sum_{i=1}^K \psi_{\phi, 1-c}^{-1}(R_{1-c}^{\phi, i}(g_i) - R_{1-c}^{\phi, i,*}), \quad \psi: \mathbb{R} \rightarrow \mathbb{R}: \text{Invertible increasing function}$$

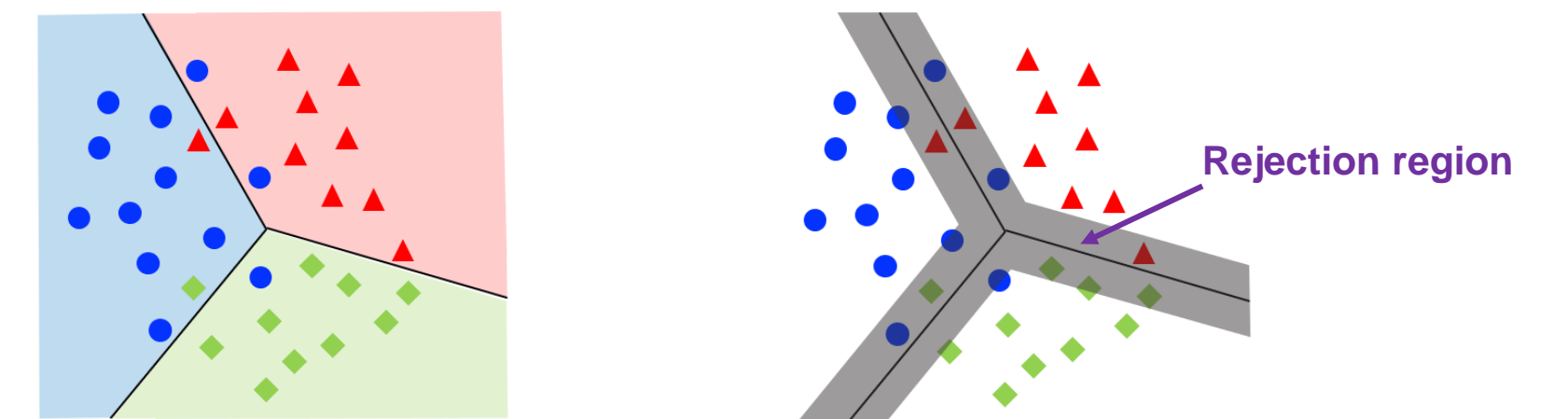
Excess cost-sensitive surrogate risk (please see our paper for more details.)

Excess 0-1-c risk is also bounded by excess cost-sensitive surrogate risk!

Connecting theory of cost-sensitive classification to classification with rejection!

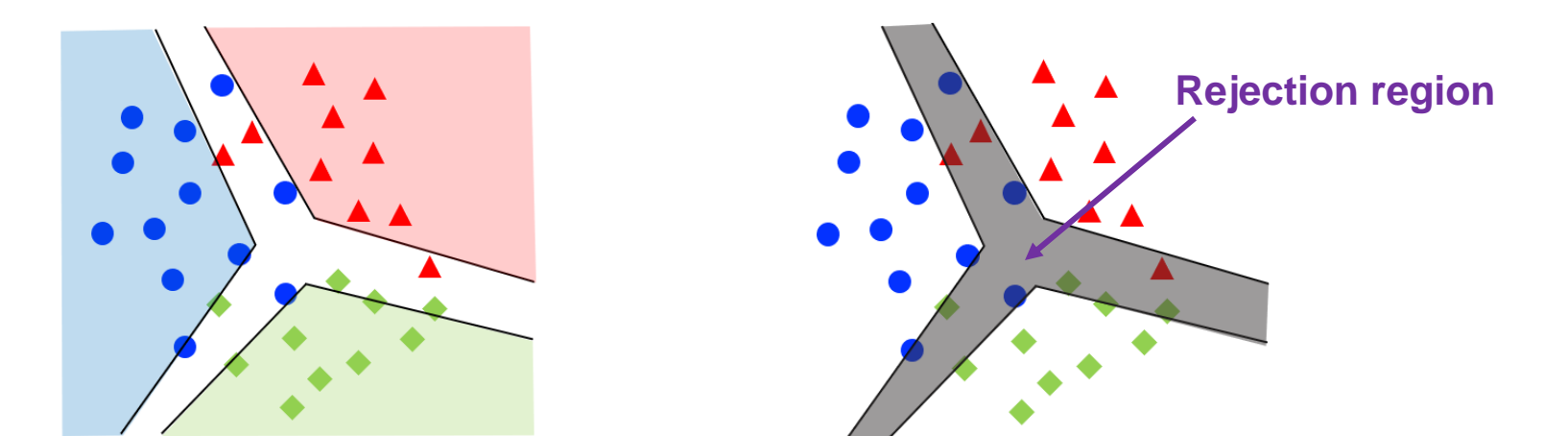
Comparison of approaches

Existing confidence-based approach



Rejection region spreads from the decision boundary.
Loss function choice is restrictive.

Proposed cost-sensitive approach



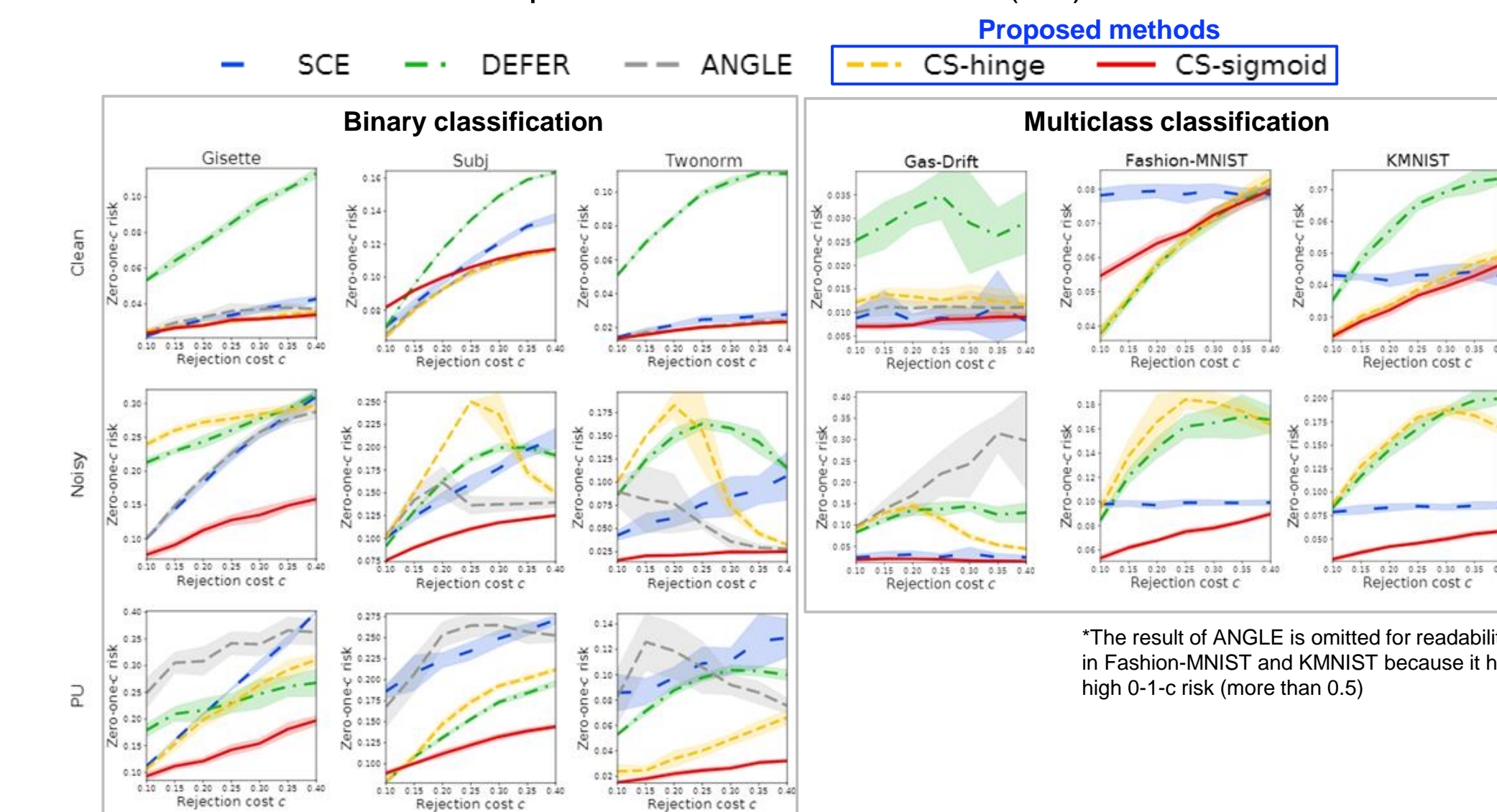
Rejection region is obtained by aggregating K -cost sensitive classifiers.
Loss requirement: **classification calibration**

Experiments

Evaluation metric: Test empirical 0-1-c risk with varying rejection cost

Baseline: Softmax cross-entropy loss with temperature scaling (SCE), DEFER (Mozannar and Sontag, 2020), ANGLE (Zhang et al., 2017)

Setting: Clean-labeled classification (Clean), Noisy-labeled classification (Noisy), Classification from positive and unlabeled data (PU)



*The result of ANGLE is omitted for readability in Fashion-MNIST and KMNIST because it has high 0-1-c risk (more than 0.5)

CS-hinge works well in classification from clean labels (Clean).

CS-sigmoid works well in classification from noisy labels (Noisy) and classification from positive and unlabeled data (PU).

*sigmoid and hinge losses are classification-calibrated but not capable of estimating $p(y|\mathbf{x})$.

References

- [1] LeCun, Y. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [2] Ni, C., Charoenphakdee, N., Honda, J., and Sugiyama, M. On the calibration of multiclass classification with rejection. *NeurIPS*, 2019.
- [3] Bartlett, P. L., and Wegkamp, M. H. Classification with a reject option using a hinge loss. *JMLR*, 2008.
- [4] Chow, C. K. On optimum recognition error and reject tradeoff. *IEEE Transactions on Information Theory*, 16(1):41–46, 1970.
- [5] Reid, M. D. and Williamson, R. C. Composite binary losses. *JMLR*, 11:2387–2422, 2010.
- [6] Scott, C. Calibrated asymmetric surrogate losses. *Electronic Journal of Statistics*, 2012.
- [7] Steinwart, I. How to compare different loss functions and their risks. *Constructive Approximation*, 2007.
- [8] Mozannar, H., and Sontag, D. Consistent estimators for learning to defer to an expert. *ICML*, 2020
- [9] Zhang, C., Wang, W., and Qiao, X. On reject and refine options in multiclass classification. *JASA*, 2017