

MOTIVATION

Selective classifier

Given a classifier $f : \mathcal{X} \rightarrow \Delta^k$, the selective classifier (f, g) at an input x is given by

$$(f, g)(x) := \begin{cases} f(x) & \text{if } g(x) \geq \tau, \\ \text{"abstain"} & \text{otherwise.} \end{cases} \quad (1)$$

where "abstain" is triggered when the *confidence scoring function* $g(x) < \tau \in \mathbb{R}$.

- **selective risk** w.r.t. $P(x, y)$ is

$$R(f, \tilde{g}) := \frac{\mathbb{E}_{P(x, y)} [\ell(f(x), y) \mathbb{I}[g(x) \geq \tau]]}{\phi(f, g)}. \quad (2)$$

- **coverage**: $\phi(f, g) = \mathbb{E}_{P(x)} [\mathbb{I}[g(x) \geq \tau]]$ represents the probability mass over the accepted samples.

Area Under the Risk Coverage curve (AURC)

The AURC [1] is typically specified as an empirical quantity from a finite sample, from which we derive the population AURC as

$$\begin{aligned} \text{AURC}_p(f) \\ = \mathbb{E}_{\tilde{x} \sim P(x)} \frac{\mathbb{E}_{(x, y) \sim P(x, y)} \ell(f(x), y) \mathbb{I}[g(x) \geq g(\tilde{x})]}{\mathbb{E}_{x' \sim P(x)} \mathbb{I}[g(x') \geq g(\tilde{x})]}. \end{aligned} \quad (3)$$

Problem

- **Finite sample limitation**: Prior works compute AURC empirically rather than at the population level, with little analysis of its statistical properties—estimator like SELE [2] can remain biased even with large sample sizes.
- **Optimization gap**: Few methods optimize AURC directly, and existing estimator do not guarantee convergence to the population AURC.

OUR INTERPRETATION

Define function $G(x)$ as the cumulative distribution function(CDF) of the CSF $g(x)$ such that

$$G(x) = \Pr(g(x') \leq g(x)) = \int \mathbb{I}[g(x') \leq g(x)] dP(x').$$

Under this definition, the population AURC in Eq. (3) is equivalent to:

$$\text{AURC}_a(f) = \int \alpha(x) \ell(f(x), y) dP(x, y) \quad (4)$$

where $\alpha(x) = -\ln(1 - G(x))$.

- **redistribution of the risk**
- $G(x)$: the **population rank percentile** based on the CSF sorted in ascending order.

PROPOSED ESTIMATORS

Estimators for $\alpha(x)$ can be achieved via Monte Carlo:

$$\hat{\alpha}_i = H_n - H_{n-r_i} \text{ and } \hat{\alpha}'_i = -\ln(1 - \frac{r_i}{n+1}). \quad (5)$$

- Both are **consistent**.
- $\hat{\alpha}_i$ upper bounds the $\hat{\alpha}'_i$, leading to

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \hat{\alpha}'_i \ell(f(x_i), y_i)}_{\widehat{\text{AURC}}'_p(f)} \leq \underbrace{\frac{1}{n} \sum_{i=1}^n \hat{\alpha}_i \ell(f(x_i), y_i)}_{\widehat{\text{AURC}}_p(f)}.$$

Proposition 1 (MSE of $\hat{\alpha}_i$ or $\hat{\alpha}'_i$) Both $\text{MSE}(\hat{\alpha}_i)$ and $\text{MSE}(\hat{\alpha}'_i)$ are asymptotically bounded by $\mathcal{O}(\frac{\beta_i}{n(1-\beta_i)+1})$.

Proposition 2 (Convergence Rate of the estimators) Assume that the loss function ℓ is square-integrable, the plug-in estimators with $\hat{\alpha}_i$ or $\hat{\alpha}'_i$ as the weight estimator, converges at a rate of $\mathcal{O}(\sqrt{\ln(n)/n})$.

EXPERIMENTS & RESULTS

Datasets. CIFAR10/100, ImageNet and a text dataset i.e Amazon Reviews.

Models. pre-trained from those datasets.

Metrics. plug-in estimators with $\hat{\alpha}$ or $\hat{\alpha}'$, the SELE score [2].

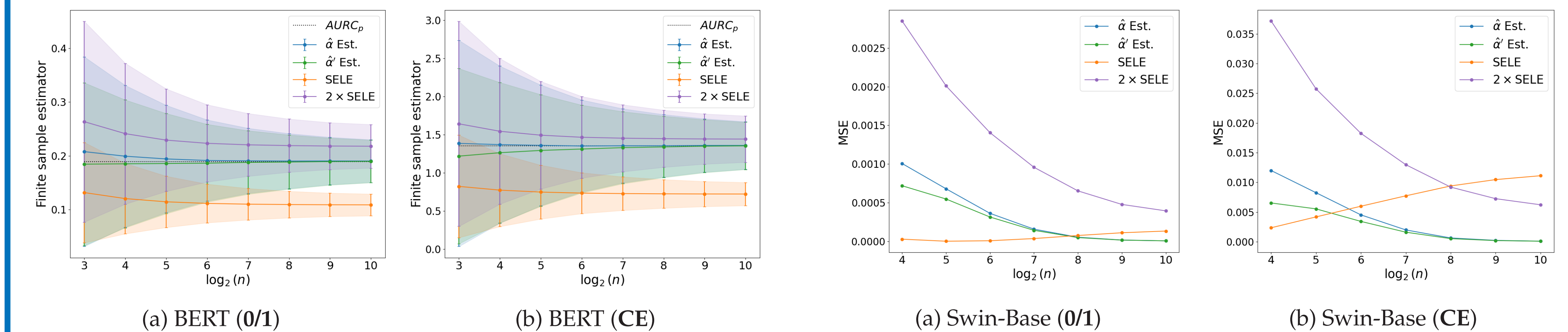


Figure 1: (Amazon) Finite-sample estimators with 0/1 or CE loss. We utilize a pre-trained model and randomly divide the test set into batch samples of size n . Subsequently, we compute the mean and standard deviation of various estimators applied to these batch samples.

Figure 2: (ImageNet) MSE of finite sample estimators with 0/1 or CE loss. For each model architecture, we calculate the MSE of the estimators using a pre-trained model on batch samples derived from the test set.

Model	CIFAR10				CIFAR100			
	CE	SELE	$\hat{\alpha}$ Est.	$\hat{\alpha}'$ Est.	CE	SELE	$\hat{\alpha}$ Est.	$\hat{\alpha}'$ Est.
ResNet18	4.967 \pm 0.038	4.470 \pm 0.030	4.473 \pm 0.030	4.471 \pm 0.030	6.648 \pm 0.021	6.577 \pm 0.011	6.532 \pm 0.012	6.533 \pm 0.014
ResNet34	6.464 \pm 0.036	5.661 \pm 0.039	5.652 \pm 0.036	5.651 \pm 0.036	6.023 \pm 0.016	5.862 \pm 0.012	5.825 \pm 0.011	5.826 \pm 0.011
ResNet50	8.318 \pm 0.002	7.892 \pm 0.046	7.921 \pm 0.047	7.918 \pm 0.049	6.225 \pm 0.009	6.043 \pm 0.015	6.007 \pm 0.008	6.007 \pm 0.009
VGG16BN	7.922 \pm 0.002	7.010 \pm 0.018	7.064 \pm 0.014	7.060 \pm 0.015	10.790 \pm 0.001	10.586 \pm 0.029	10.559 \pm 0.029	10.560 \pm 0.030
VGG19BN	9.813 \pm 0.192	8.475 \pm 0.061	8.528 \pm 0.059	8.524 \pm 0.059	10.633 \pm 0.001	10.421 \pm 0.026	10.393 \pm 0.025	10.391 \pm 0.024
WideResNet28x10	4.137 \pm 0.046	3.867 \pm 0.049	3.864 \pm 0.049	3.863 \pm 0.049	5.912 \pm 0.652	5.607 \pm 0.707	5.836 \pm 0.652	5.607 \pm 0.707

Table 1: Summary of population AURC_p (mean \pm std, scaled by 10^{-2}) on the test set for models fine-tuned with various loss functions. Each entry aggregates results over five seeds using the same pre-trained model.

CONCLUSION

- We extend empirical AURC to a true population quantity and show it admits a reweighted-risk interpretation.
- We propose two plug-in estimators via Monte Carlo method and show their bias, MSE, and an $\mathcal{O}(\sqrt{\ln n/n})$ convergence rate.
- Experiments demonstrate these consistent estimators not only outperform SELE in terms of estimation but also serve as effective objectives for directly fine-tuning networks to minimize AURC.

REFERENCES

- [1] Yonatan Geifman, Guy Uziel, and Ran El-Yaniv. Bias-reduced uncertainty estimation for deep neural classifiers. In *International Conference on Learning Representations*, 2019.
- [2] Vojtech Franc, Daniel Prusa, and Vaclav Voracek. Optimal strategies for reject option classifiers. *Journal of Machine Learning Research*, 24(11):1–49, 2023.

CONTACT INFORMATION

Paper <https://arxiv.org/pdf/2410.15361>
Email han.zhou@esat.kuleuven.be