

Deep Neural Network Benchmarks for Selective Classification

Andrea Pugnana

*Scuola Normale Superiore, University of Pisa, ISTI-CNR
Pisa, Italy*

ANDREA.PUGNANA@DI.UNIPI.IT

Lorenzo Perini

*KU Leuven
Leuven, Belgium*

LORENZO.PERINI@KULEUVEN.BE

Jesse Davis

*KU Leuven
Leuven, Belgium*

JESSE.DAVIS@KULEUVEN.BE

Salvatore Ruggieri

*University of Pisa
Pisa, Italy*

SALVATORE.RUGGIERI@UNIPI.IT

<https://openreview.net/forum?id=xDPzHbtAEs>

Editor: Mykola Pechenizkiy

Abstract

With the increasing deployment of machine learning models in many socially-sensitive tasks, there is a growing demand for reliable and trustworthy predictions. One way to accomplish these requirements is to allow a model to abstain from making a prediction when there is a high risk of making an error. This requires adding a selection mechanism to the model, which *selects* those examples for which the model will provide a prediction. The selective classification framework aims to design a mechanism that balances the fraction of rejected predictions (i.e., the proportion of examples for which the model does not make a prediction) versus the improvement in predictive performance on the selected predictions. Multiple selective classification frameworks exist, most of which rely on deep neural network architectures. However, the empirical evaluation of the existing approaches is still limited to partial comparisons among methods and settings, providing practitioners with little insight into their relative merits. We fill this gap by benchmarking 18 baselines on a diverse set of 44 datasets that includes both image and tabular data. Moreover, there is a mix of binary and multiclass tasks. We evaluate these approaches using several criteria, including selective error rate, empirical coverage, distribution of rejected instance's classes, and performance on out-of-distribution instances. The results indicate that there is not a single clear winner among the surveyed baselines, and the best method depends on the users' objectives.

1 Introduction

Artificial Intelligence (AI) systems are increasingly being deployed to support or even automate decision-making. Ensuring the trustworthiness of AI systems is crucial in many applications (Kaur et al., 2023), and is one of the main goals of the recent European AI Act (European Commission, 2021). More precisely, “[h]igh-risk AI systems shall be designed and developed in such a way that they achieve, in the light of their intended purpose, an appropriate level of accuracy [and] robustness”.

High-risk AI systems pertain to socially sensitive domains, such as: healthcare, where predictions might be used to determine treatments (Craig et al., 2023); justice, where predictions can evaluate the risk of recidivism (Berk et al., 2021); hiring, where predictions can determine rankings of candidates or explain their turnover intention (Fabris et al., 2023; Lazzari et al., 2022); and credit scoring, where predictions can be used to estimate the probability of repaying a debt (Dastile et al., 2020).

In all such high-risk contexts, we aim to reduce the number of mistakes made by AI systems because their mistakes can have critical consequences. For example, consider a bank that uses a Machine Learning (ML) model to score the credit risk of loan applications. In such a setting, a misprediction could either translate into a money loss for the bank or an unjust denial of credit to the applicant.

One potential way to improve the trustworthiness of a model is to allow it to abstain from making a prediction when there is a high chance of making an error (Chow, 1970). Such a strategy is inherent in human reasoning when facing an unknown phenomenon. For example, human bankers who are unsure about a specific loan application do not (have to) provide an answer as soon as they are asked. Indeed, they may require additional financial documents to verify the loan’s feasibility or ask for an external expert consultation. This approach aims to minimize the risk of an incorrect evaluation.

Likewise, allowing ML models to predict only when confident enough helps mitigate the risk of incorrect predictions (Pugnana, 2023). On the one hand, including a reject option results in the ML model having better performance when it does provide a prediction because it is only offering predictions in those cases where it is highly likely to be correct. On the other hand, rejected instances can be dealt with in other ways. For example, human experts can be involved in the loop to oversee difficult instances, e.g., a banker can oversee difficult-to-evaluate loan applications. Alternatively, the prediction task can be deferred to more complex ML models, possibly using additional and costly-to-compute features.

Selective Classification (SC) (El-Yaniv and Wiener, 2010) is one well-known framework that allows a model not always to offer a prediction. Intuitively, this framework imbues a model with a mechanism that *selects* whether a prediction is made on a per-example basis. The goal is to navigate the tradeoff between the proportion of examples for which a prediction is made (i.e., the model’s *coverage*) and the performance improvement on the selected examples (i.e., the ones for which a prediction is made) that arises from focusing only on those cases where the model has a small chance of making a misprediction. Typically, this is done by maximizing the performance on the selected examples given a target coverage. Given the appeal of SC, there are wide range of approaches for this problem setting (Geifman and El-Yaniv, 2017, 2019; Liu et al., 2019; Huang et al., 2020; Gangrade et al., 2021; Pugnana

and Ruggieri, 2023a,b; Feng et al., 2023). The primary emphasis is on implementing SC in the context Deep Neural Networks (DNN) models.

Unfortunately, we lack insights into the relative merits of existing SC approaches for DNNs because existing empirical evaluations in the literature suffer from several shortcomings. First, they always involve ≤ 10 datasets, and primarily consider only image data. Second, only a handful of approaches (never more than seven) are compared. Third, most studies mainly focus on comparing approaches based on single metric: their predictive accuracy on the selected examples. However, there are other relevant performance characteristics of SC methods such as whether their coverage constraint holds, whether they disproportionately reject instances from one class, or how they behave on unseen data.

Our goal is to fill this gap by performing the first comprehensive benchmarking of SC methods for DNN architectures. Specifically, our evaluation goes substantially beyond existing studies by:

1. Including 18 SC methods;
2. Evaluating the considered methods on 44 datasets that include both image and tabular data; and
3. Considering five different aspects of SC models' performance.

Our results suggest that the choice of the baseline depends on the performance criterion to be prioritized. In fact, most methods perform with no statistically significant difference across the different tasks. To summarize, the main contributions of this paper are that we:

- (i) briefly survey the state-of-the-art methods in SC;
- (ii) provide the widest experimental evaluation of SC methods in terms of baselines, datasets and tasks;
- (iii) point out the limitations of compared methods, which highlights potential avenues for future research directions; and
- (iv) release a public repository with all software code and datasets for reproducing the baseline algorithms and the experiments.¹

2 Background

Let \mathcal{X} be an d -dimensional input space, $\mathcal{Y} = \{1, \dots, m\}$ be the target space and $P(\mathbf{X}, Y)$ be the probability distribution over $\mathcal{X} \times \mathcal{Y}$. Given a hypothesis space \mathcal{H} of functions that map \mathcal{X} to \mathcal{Y} (called models or classifiers), the goal of a learning algorithm is to find the hypothesis $h \in \mathcal{H}$ that minimizes the *risk*:

$$R(h) = \mathbb{E}[l(h(\mathbf{X}), Y)] \tag{1}$$

where $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a user-specified loss function. Because $P(\mathbf{X}, Y)$ is generally unknown, it is typically assumed that we have access to an i.i.d. sample $\mathcal{T}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ that can

1. The code is available at github.com/andrepugni/ESC/.

be used to learn a classifier $\hat{h}(\cdot)$, such that:

$$\hat{h} \in \arg \min_{h \in \mathcal{H}} \hat{R}(h, \mathcal{T}_n) \quad (2)$$

where $\hat{R}(h, \mathcal{T}_n) = 1/|\mathcal{T}_n| \sum_{(\mathbf{x}, y) \in \mathcal{T}_n} l(h(\mathbf{x}), y)$ is the *empirical risk* over the sample \mathcal{T}_n .

Because the learned model is prone to making mistakes, one can extend the canonical setting to include a selection mechanism that allows the model to refrain from offering a prediction for those instances likely to be misclassified.

Formally, a *selective classifier* is a pair (h, g) where h is a standard classifier and $g : \mathcal{X} \rightarrow \{0, 1\}$ is a *selection function* that determines whether h 's prediction is provided or the model abstains (or rejects):

$$(h, g)(\mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{if } g(\mathbf{x}) = 1 \\ \text{abstain} & \text{otherwise} \end{cases} \quad (3)$$

In practice, rather than directly learning the selection function in Eq. 3, one approximates it by (1) learning a *confidence function*² $k_h : \mathcal{X} \rightarrow [0, 1]$ (sometimes called soft selection (Geifman and El-Yaniv, 2017)) that measures how likely it is that the predictor h is correct, and (2) setting a threshold $\tau \in [0, 1]$ that defines the minimum confidence for providing a prediction. A low confidence value indicates that the model is likely to make a misprediction for the instance and therefore it should abstain, which yields the following selection function:

$$g(\mathbf{x}) = \mathbb{1}(k_h(\mathbf{x}) > \tau) \quad (4)$$

To prevent the selective classifier from abstaining on too many (test) instances, SC methods also consider the *coverage* metric, which is defined as

$$\phi(g) = \mathbb{E}[g(\mathbf{X})]. \quad (5)$$

The coverage computes the expected proportion of instances for which the model would make a prediction. These non-rejected instances are commonly referred to as either *accepted* or *selected*, and we will use these terms interchangeably. We will refer to the *rejection rate* as the complement of the coverage, i.e., $1 - \phi(g)$ (Perini and Davis, 2023). Another core measure of the SC framework is the risk over the accepted region, commonly called the *selective risk* which is defined as:

$$R(h, g) = \frac{\mathbb{E}[l(h(\mathbf{X}), Y)g(\mathbf{X})]}{\phi(g)} = \mathbb{E}[l(h(\mathbf{X}), Y)|g(\mathbf{X}) = 1] \quad (6)$$

A widely adopted instance of the selective risk is the *selective error rate*, which corresponds to the selective risk for the 0-1 loss $l(h(\mathbf{X}), Y) = \mathbb{1}\{h(\mathbf{X}) \neq Y\}$.

Coverage and risk are estimated over a given test set \mathcal{T}_{test} as follows. The *empirical risk* over the set of accepted instances is defined as:

$$\hat{R}(h, g, \mathcal{T}_{test}) = \frac{1}{|\mathcal{T}_{test}| \cdot \hat{\phi}(g, \mathcal{T}_{test})} \sum_{(\mathbf{x}, y) \in \mathcal{T}_{test}} l(h(\mathbf{x}), y) \cdot g(\mathbf{x}) \quad (7)$$

2. A good confidence function k_h should rank instances based on descending loss, i.e., if $k_h(\mathbf{x}_i) \leq k_h(\mathbf{x}_j)$ then $l(h(\mathbf{x}_i), y_i) \geq l(h(\mathbf{x}_j), y_j)$.

where $\hat{\phi}(g, \mathcal{T}_{test}) = 1/|\mathcal{T}_{test}| \sum_{(\mathbf{x}, y) \in \mathcal{T}_{test}} g(\mathbf{x})$ is the *empirical coverage* over the test set. Observe that $\hat{R}(h, g, \mathcal{T}_{test}) = \hat{R}(h, \mathcal{T}_{test}^g)$, where $\mathcal{T}_{test}^g = \{(\mathbf{x}, y) \in \mathcal{T}_{test} \mid g(\mathbf{x}) = 1\}$, i.e., the empirical risk of a selective classifier boils down to the empirical risk of the classifier over the set of accepted instances. The inherent trade-off between coverage and risk can be summarized by a *risk-coverage curve* (El-Yaniv and Wiener, 2010). Moreover, this trade-off allows framing the SC task according to two different formulations: the bounded improvement model and the bounded abstention model (Franc et al., 2023). In the bounded improvement model, the problem is formulated by fixing an upper bound r - the *target risk* - for the selective risk and then looking for a selective classifier that maximizes coverage (Geifman and El-Yaniv, 2017).

Problem 1 (Bounded-improvement model) *Given a target risk r , an optimal selective classifier (h, g) is a solution to:*

$$\max_{\theta, \psi} \phi(g_\psi) \quad s.t. \quad R(h_\theta, g_\psi) \leq r \quad (8)$$

Conversely, in the bounded-abstention model, we fix a lower bound c for coverage (called *target coverage*) and then look for a selective classifier that minimizes the selective risk (Geifman and El-Yaniv, 2019).

Problem 2 (Bounded-abstention model) *Given a target coverage c , an optimal selective classifier (h, g) is a solution to:*

$$\min_{\theta, \psi} R(h_\theta, g_\psi) \quad s.t. \quad \phi(g_\psi) \geq c \quad (9)$$

We call *coverage-calibration* the post-training procedure of estimating the threshold τ in (4) for the target coverage c specified in Problem 2. This is generally done by estimating the $(1 - c) \cdot 100$ -th percentile of the confidence function k_h over a held-out calibration set \mathcal{T}_{cal} .

3 Baselines

There are multiple ways to devise abstaining classifiers. We restrict our attention to DNN approaches aiming to solve the bounded-abstention problem (Eq. 9). We present and categorize a few baselines according to their definition of the confidence function, extending the work of Feng et al. (2023). We distinguish among three categories of methods: **Learn-to-Abstain**, **Learn-to-Select** and **Score-based**.

3.1 Learn-to-Abstain Methods

Learn-to-Abstain methods tackle the selective classification task by adding a new class label $(m + 1)$ representing abstention to the classification problem. While there are no actual instances belonging to this class, these approaches design loss functions to enable the classifier to assign a positive score $s_{m+1}(\mathbf{x})$ to ambiguous instances. This score serves as a confidence function, i.e., $k_h(\mathbf{x}) = 1 - s_{m+1}(\mathbf{x})$ (Feng et al., 2023). In Figure 1, we provide an example of a canonical Learn-to-Abstain architecture.

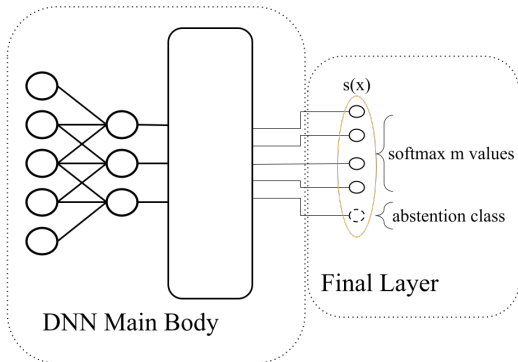


Figure 1: A generic Learn-to-Abstain architecture

The first method to take this approach was **DG** (Liu et al., 2019). It uses a reward hyperparameter o for the class $m+1$ to set how often the classifier should abstain. Formally, DG trains a neural network minimizing the following loss:

$$\mathcal{L}_{\text{DG}} = \mathbb{E}_{P(\mathbf{x}, Y)} \left[\log(s_y(\mathbf{x})) + \frac{1}{o} s_{m+1}(\mathbf{x}) \right], \quad (10)$$

where $s_y(\mathbf{x})$ and $s_{m+1}(\mathbf{x})$ are the neural network softmax values, respectively, over the true class $Y = y$ and $m+1$ (abstention). Intuitively, a higher o encourages the network to be confident in its prediction, and a low o makes it less confident and more likely to abstain. However, DG does not have any explicit way to guide abstention towards more difficult examples during training, as the reward o remains fixed for the whole training procedure.

To overcome this limitation, Self-Adaptive Training (**SAT**) (Huang et al., 2020) trains the selective classifier through a convex combination of predictions and true labels. This combination is dynamically adapted during the training process to identify those instances that are more difficult to correctly classify and, hence are good candidates for abstention. More precisely, for the first E_s (user-defined) epochs, the training target - $\mathbf{t} \in [0, 1]^m$ - is equal to the one-hot encoded true label vector \mathbf{y} . Afterwards, it becomes the convex combination of (probabilistic) predictions and true labels, namely $\mathbf{t} = \gamma \mathbf{t} + (1 - \gamma) \mathbf{s}(\mathbf{x})$, with $\mathbf{s}(\mathbf{x})$ representing the neural network softmax values and γ the weight of the convex combination. The final selective classifier is then optimized by minimizing the loss function:

$$\mathcal{L}_{\text{SAT}} = -\mathbb{E}_{P(\mathbf{x}, Y)} \left[\mathbf{t}' \log(\mathbf{s}(\mathbf{x})) + (1 - t_y) \log s_{m+1}(\mathbf{x}) \right], \quad (11)$$

where t_y is the value of vector \mathbf{t} corresponding to the index of true value y and $s_{m+1}(\mathbf{x})$ represents the softmax value for the abstention class.³ Both DG and SAT add an extra softmax value to the neural network output to identify difficult-to-predict instances. However,

3. For instance, if $y = 1$ and $m = 2$, then $\mathbf{t}' = [0, 1]$ and $t_y = 1$ when epoch is below E_s . Intuitively, the first term is the cross-entropy loss between the classifier and the adaptive training target, which allows learning a good multi-class classifier. The second term serves as a confidence function, identifying uncertain samples in the dataset. The balance between these terms is controlled by the value of t_y , which determines whether the classifier learns to abstain or make accurate predictions.

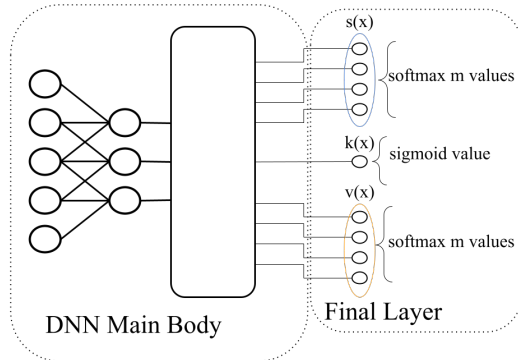


Figure 2: An example of SELNET, a Learn-to-Select architecture.

Feng et al. (2023) argue that incorporating this extra class in the training loss could result in overfitting on examples that are easier to classify. To mitigate this, **SAT+EM** (Feng et al., 2023) adds an average entropy term $\mathcal{E}(\mathbf{s}(\mathbf{x}))$ to SAT’s loss :

$$\mathcal{L}_{\text{SAT+EM}} = \mathcal{L}_{\text{SAT}} + \beta \mathcal{E}(\mathbf{s}(\mathbf{x})) \quad (12)$$

where $\mathbf{s}(\mathbf{x})$ represents the neural network of m softmax values, and β is a hyperparameter that measures the impact of the entropy term. All the learn-to-abstain methods are calibrated for the target coverage c using a calibration set (as discussed in Section 2).

3.2 Learn-to-Select Methods

Like Learn-to-Abstain methods, Learn-to-Select methods simultaneously learn the classifier and its specific confidence function. However, in this setting, the confidence function does not rely on an additional abstention class but aims at achieving a specific target coverage c . This procedure ensures that the classifier’s parameters are optimized to correctly predict instances less likely to be rejected.

The main architecture belonging to this class is SelectiveNet (**SELNET**) (Geifman and El-Yaniv, 2019). Given a target coverage c , SELNET jointly trains the final classifier and the confidence function to maximize the performance over the $100 \cdot c\%$ most confident instances. SELNET’s architecture has four main components, each with a different purpose, as depicted in Figure 2: the main body, the predictive head s , the selective head k , and the auxiliary head v . The main body consists of deep layers shared by all three heads: any deep-learning architecture can be used in this part (e.g., convolutional layers, linear layers, recurrent layers, etc.). The predictive head $s(\mathbf{x})$, consisting of a final linear layer with softmax, is used to make the classifier prediction. The selective head $k(\mathbf{x})$ outputs a confidence function using a linear layer with a final sigmoid activation. The auxiliary head $v(\mathbf{x})$ replicates the structure of the predictive head and mitigates the risk of overfitting on the accepted instances. Given the target coverage c , SELNET is trained using the following

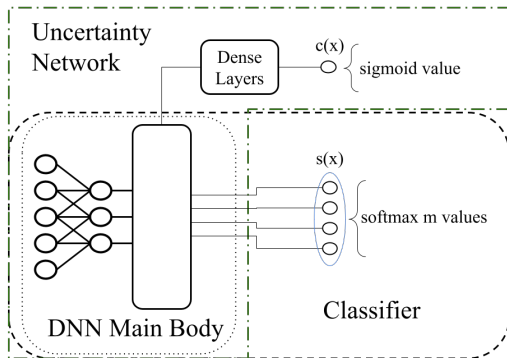


Figure 3: An example of CONFIDNET, a score-based architecture.

loss function:

$$\mathcal{L}_{\text{SELNET}} = \alpha \left(\frac{\mathbb{E}_{P(\mathbf{x}, Y)} [l(s(\mathbf{x}), y)k(\mathbf{x})]}{\phi(k)} + \lambda(\max(0, c - \phi(k)))^2 \right) + (1-\alpha)\mathbb{E}_{P(\mathbf{x}, Y)} [l(v(\mathbf{x}), y)] \quad (13)$$

where $l(s(\mathbf{x}), y)$ is the cross-entropy loss for the predictive head $s(\mathbf{x})$; $\phi(k)$ is the coverage obtained by selective head k ; $l(v(\mathbf{x}), y)$ is the cross entropy loss for auxiliary head prediction $v(\mathbf{x})$; α is a hyperparameter to control the relative importance between the losses for the predictive and the auxiliary head; and λ is a penalization term for coverage violations.

For the sake of completeness, following the same reasoning as for SAT+EM (Feng et al., 2023), we include **SELNET+EM** in the comparison. This approach adapts the SELNET objective function to contain an additional entropy term.

3.3 Score-based Methods

Score-based methods compute and set a threshold on a confidence function - as formalized in Eq. 4 - that is based on the classifier’s output. Conceptually, this means that predictions are only made for the test examples for which the model is most confident. Because this can be viewed as a post-hoc approach, this confers the advantage of being applicable to already trained models.

The most popular technique is the Softmax Response **SR** (Geifman and El-Yaniv, 2017), which defines the confidence function as the maximum value of a final softmax layer, i.e., $k_{\text{SR}}(x) = \max_{y \in \mathcal{Y}} s_y(\mathbf{x})$. Given a coverage c , SR sets the selection threshold τ using the calibration procedure explained in Section 2.

Since the SR principle is very general, it can be applied to any method that provides scores for the classes (Franc et al., 2023). For example, Feng et al. (2023) propose to improve learn-to-abstain and learn-to-select methods by replacing their confidence function with the SR confidence. In particular, three novel methods are presented, i.e., **SAT+SR**, **SAT+EM+SR**, **SELNET+SR**, which are trained using \mathcal{L}_{SAT} , $\mathcal{L}_{\text{SAT+EM}}$ and $\mathcal{L}_{\text{SELNET}}$ respectively. For the sake of completeness, we include also **SELNET+EM+SR** in the comparison, i.e., a network trained with the SELNET+EM’s loss and using the SR selection strategy.

Another score-based popular option is using ensembles of neural networks. For example, Lakshminarayanan et al. (2017) train multiple networks with different initialization and build a selective classifier by computing the entropy of the (multiple) network outputs (**ENS**). The intuition is that more disagreement among the outputs indicates that the ensemble is uncertain about its prediction, and hence rejection is appropriate. However, despite the advantages of using ensembles in terms of performance, relating a dispersion measure to the correctness of predictions is not straightforward. Hence, the authors also propose using the average softmax response (i.e., $k_{\text{ENS}}(\mathbf{x}) = 1/J \sum_{j=1}^J k_{\text{SR},j}(\mathbf{x})$, where J is the number of networks in the ensemble) as a confidence measure. We will refer to this baseline as **ENS+SR**. A theoretical analysis of the advantages of using ENS+SR can be found in Ding et al. (2023).

The main concern with using $k_{\text{SR}}(\mathbf{x})$ as a confidence measure is that may provide high values both for mistakes and correct predictions, making them indistinguishable. On the other hand, when the model misclassifies an example, the score $s_y(\mathbf{x})$ associated with the true class probability $P(Y = y | \mathbf{X} = \mathbf{x})$ should be low, making it a viable option to perform selective classification. However, one cannot access true labels at test time, making it impossible to use $s_y(\mathbf{x})$ directly. Corbière et al. (2019) address this concern by estimating $s_y(\mathbf{x})$ with a two-step procedure called **CONFIDNET** (as depicted in Figure 3). First, they estimate $s_y(\mathbf{x})$ by training a neural network classifier. Next, they build a second (uncertainty) network on top of the classifier: the main body is kept unchanged, while the final part of the original classifier is replaced with a series of dense layers. This uncertainty network is then trained considering the following loss function:

$$\mathcal{L}_{\text{CONFIDNET}} = \mathbb{E}_{P(\mathbf{X}, Y)}[(c(\mathbf{x}) - s_y(\mathbf{x}))^2] \quad (14)$$

with $c(\mathbf{x})$ referring to the final output of the uncertainty network. Intuitively, $c(\mathbf{x})$ should mimic $s_y(\mathbf{x})$ and can be used as a confidence function: the higher $c(\mathbf{x})$, the higher the chance the classifier is right.

Franc et al. (2023) also use a classifier and an uncertainty estimator. They propose two different approaches, named **REG** and **SELE**. Both of them learn the classifier on half of the training data and use the other half to directly estimate where the classifier is more likely to make mistakes. In particular, these two methods focus on learning an uncertainty score f , which mirrors the confidence function k : the higher f is, the higher the likelihood of making mistakes (thus, abstention is preferable in the latter). Neither SELE nor REG are tied to specific neural network architectures, i.e., they are model-agnostic and can be adapted to other learning models. REG poses the problem of learning the uncertainty score as a regression problem, where given a set of hypotheses \mathcal{F} , the uncertainty score $f \in \mathcal{F}$ minimizes the following:

$$\mathcal{L}_{\text{REG}} = \mathbb{E}_{P(\mathbf{X}, Y)}[(l(h(\mathbf{x}, y)) - f(\mathbf{x}))^2] \quad (15)$$

Intuitively, the higher the value of f , the higher the loss. Hence, abstention should be preferred. On the other hand, given a hypothesis space \mathcal{F} , SELE considers a surrogate loss of the risk coverage curve, i.e.,

$$\mathcal{L}_{\text{SELE}} = \mathbb{E}_{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \sim P(\mathbf{X}, Y)}[l(h(\mathbf{x}_1, y_1)) \log(1 + \exp(f(\mathbf{x}_2) - f(\mathbf{x}_1)))] \quad (16)$$

and then learns the uncertainty score $f \in \mathcal{F}$ by minimizing $\mathcal{L}_{\text{SELE}}$.

The approaches presented so far require a held-out calibration dataset. Unfortunately, for problems where only little data is available, reducing the amount of training data may deteriorate the classifier’s performance. Moreover, splitting the data into a dataset for training and a dataset for calibration may introduce randomness effects on both the classifier and the selection function. **SCROSS** (Pugnana and Ruggieri, 2023a) is a model-agnostic approach that overcomes the need for a calibration set by employing a cross-validation strategy that follows three steps. First, it splits the available data into K folds. Second, it trains a classifier over $K - 1$ folds and predicts the SR confidence values over the remaining K -th fold. Finally, it stacks the predicted confidence values altogether. This approach approximates the confidence over the full dataset. Then, SCROSS uses SR’s approach to threshold such confidence values.

Moreover, in high-risk scenarios where SC is sought, such as healthcare and finance, we often deal with imbalanced (binary) classes (He and Garcia, 2009). A common metric to evaluate the performance of classifiers in such contexts is the Area Under the ROC Curve (AUC) (Yang and Ying, 2023), which measures the classifier’s ability to rank instances from minority and majority classes correctly. Pugnana and Ruggieri (2023b) provide a theoretical condition - two bounds over the minority class score - that guarantees not to worsen AUC once we allow for abstention. The selection function is implemented by (1) estimating these lower and upper bounds for the minority class score, and (2) rejecting instances with minority class scores between the two (estimated) bounds. To implement such a strategy, the authors devise two algorithms, i.e., **PLUGINAUC** and **AUCROSS**. The difference between the two methods lies in how their selection strategy is calibrated: PLUGINAUC adopts a held-out approach to calibrate the bounds, while AUCROSS uses a cross-fitting approach similar to SCROSS.

4 Research Questions

This paper intends to evaluate the relative strength of the baselines introduced in Section 3 with respect to the following research questions:

- Q1:** Are there significant differences across baselines and scenarios regarding selective error rate?
- Q2:** Are there significant differences across baselines and scenarios regarding violations of the target coverage?
- Q3:** How are the methods’ rejection rates distributed among the classes?
- Q4:** How do the methods behave when flipping the learning task to maximise the coverage under constraints on the error rate?
- Q5:** How do the methods react to out-of-distribution test examples?

We differentiate from previous works in several respects:

- Regarding **Q1**, our study goes beyond existing ones in two important ways. First, prior evaluations involving SC methods were performed using less than seven baselines

and less than ten datasets whereas we consider 18 methods and 44 datasets. Second, prior benchmarks largely focused on image data whereas our benchmark also include tabular data.

- Concerning **Q2**, only a few works investigate coverage violations, i.e., Geifman and El-Yaniv (2019); Pugnana and Ruggieri (2023a,b). As in **Q1**, this was done on a much smaller scale: for example, Geifman and El-Yaniv (2019) considered only a single image dataset, while Pugnana and Ruggieri (2023a) and Pugnana and Ruggieri (2023b) considered eight and nine binary datasets respectively;
- Only the work by Pugnana and Ruggieri (2023b) addresses **Q3** and highlights that the rejection rate is biased against the minority class. However, they considered nine binary datasets and only six baselines;
- We are the first to empirically evaluate **Q4** and assess performances when switching from minimizing selective risk to maximizing coverage on a large and diverse set of data and settings;
- We are the first to evaluate **Q5** and evaluate how SC methods perform when dealing with shifts in the feature space.

5 Experimental Evaluation

5.1 Experimental Setting

Datasets and Baselines. We run experiments on 44 benchmark datasets from real-life scenarios, such as finance and healthcare (Yang et al., 2023). Among these, 20 are image data and 24 are tabular data. Moreover, 13 of these datasets were previously used in testing (at least one) the baselines in their original paper. Details are provided in Tables A1-A2 of the Appendix A.1. We compare a total of 18 baseline methods (presented in Section 3) representing the state-of-the-art SC methods: DG, SAT, SAT+EM (*learn-to-abstain*); SELNET, SELNET+EM (*learn-to-select*); SR, SAT+SR, SAT+EM+SR, SELNET+SR, SELNET+EM+SR, ENS, ENS+SR, CONFIDNET, REG, SELE, SCROSS, PLUG-INAUC, AUCROSS (*score-based*).

Hyperparameters. The baselines share the same neural-network architecture. For image data, we use either a Resnet34 architecture (He et al., 2016) or the one specified in the original paper. For tabular data, since neural networks are not state-of-the-art methods, we use the architectures proposed by Gorishniy et al. (2021); Grinsztajn et al. (2022), which revised DNN models for tabular data. Overall, we consider two sets of hyperparameters: network-specific (e.g., hidden layers, learning rate), and loss-specific (e.g., β for SAT+EM). All networks are trained for 300 epochs. We optimize the hyperparameters using *optuna* (Akiba et al., 2019), a framework for multi-objective Bayesian optimization, with the following inputs: coverage violation and cross-entropy loss as target metrics, *BoTorch* as sampler (Balandat et al., 2020), 10 initial independent trials out of 20 total trials. Among the 20 trials, we select the configuration that (1) has the highest accuracy on the validation set and (2) reaches the target coverage (± 0.05). Moreover, some baselines require the target coverage c to be known at training time (e.g., SELNET). For the

sake of reducing the computational cost⁴, we optimize their hyperparameters using only three values $c \in \{.99, .85, .70\}$ and fix the best-performing architecture for all target coverages. Moreover, SCROSS, AUCROSS, ENS, ENS+SR and PLUGINAUC use the same optimal hyperparameters found for SR as they share the same training loss. Similarly, SELNET+SR, SELNET+EM+SR, SAT+SR and SAT+EM+SR employ the same optimal configuration as, respectively, SELNET, SELNET+EM, SAT and SAT+EM. We detail the parameter choices in Appendix A.2.

Experimental setup. For each combination of datasets and baselines, we run the following experiment: (i) we randomly split the available data into training, calibration, validation, and test sets using the proportion 60/10/10/20%, (ii) we consider the following 7 target coverages $c \in \{.7, .75, .8, .85, .9, .95, .99\}$, (iii) we tune the baseline’s hyperparameters using training, calibration, and validation sets as described in the previous paragraph, (iv) we use such optimal hyperparameters to train the baseline on the training set and calibrate the confidence function on the calibration set, (v) we draw 100 bootstrapped datasets from the test set (see (Rajkomar et al., 2018)) with the same size at the test set, and, finally, (vi) we compute the empirical selective error rate⁵ $\widehat{Err}(h, g, \mathcal{T}_{test})$, the empirical coverage $\widehat{\phi}(g, \mathcal{T}_{test})$, and, for binary datasets, the class distribution over the accepted instances for each of the 100 bootstrapped datasets \mathcal{T}_{test} . For each evaluation metric, we compute its mean and standard deviation over the 100 bootstrap runs. In reporting results, we distinguish between binary and multi-class (i.e., > 2 classes) problems because PLUGINAUC and AUCROSS are specific for binary classification.

Regarding computational resources, we split the workload over three machines: (1) a 25 nodes cluster equipped with 2×16 -core @ 2.7 GHz (3.3 GHz Turbo) POWER9 Processor and 4 NVIDIA Tesla V100 each, OS RedHatEnterprise Linux release 8.4; (2) a 96 cores machine with Intel(R) Xeon(R) Gold 6342 CPU @ 2.80GHz and two NVIDIA RTX A6000, OS Ubuntu 20.04.4; (3) a 128 cores machine with AMD EPYC 7502 32-Core Processor and four NVIDIA RTX A5000, OS Ubuntu 20.04.6.

5.2 Experimental Results

We report here the main experimental results w.r.t. the research questions Q1–Q5. Additional results are reported in the Appendix B.

Q1. Comparing the error rates. We introduce a normalized version of the empirical selective error rate, called *relative error rate*:

$$RelErr(h, g, \mathcal{T}_{test}) = \frac{\widehat{Err}(h, g, \mathcal{T}_{test})}{\widehat{Err}(h_{maj}, g, \mathcal{T}_{test})}, \quad (17)$$

where $\widehat{Err}(h_{maj}, g, \mathcal{T}_{test})$ is the empirical selective error rate obtained by always predicting the majority class in the training set. This normalization accounts for variability in task

4. Tuning the networks is computationally expensive, requiring more than 15 days on some large datasets, such as `food101`.

5. The empirical selective error rate is the empirical risk (7) w.r.t. the 0-1 loss. Almost all of the baselines are optimized for such a metric, except PLUGINAUC and AUCROSS that are designed for increasing AUC.

prediction difficulty. Intuitively, the closer the relative error rate to 0 the better. Values close to 1 denote selective error rates similar to the ones of a majority classifier.

Figure 4 reports the mean relative error rates for the top two and the worst two⁶ baselines. We limit the number of reported baselines for clarity. Tables with detailed results at the dataset level are reported in the Appendix B.3.

For binary data, the best-performing methods are ENS+SR and SR. ENS+SR’s relative error rate is $\approx .485$ at $c = .99$, decreasing to $\approx .365$ at $c = .70$. SR ranges from $\approx .488$ at $c = .99$ to $\approx .363$ at $c = .70$. The worst baselines are DG and REG, with relative error rates of $\approx .615$ and $.544$ at $c = .99$ respectively, up to $\approx .564$ and $\approx .529$ at $c = .70$.

Also for multiclass data, ENS+SR and SR achieve the best results. The relative error rate ENS+SR ranges from $\approx .182$ at $c = .99$ to $\approx .117$ at $c = .70$, while SR starts from $\approx .197$ at $c = .99$, and decreases down to $\approx .127$ for $c = .70$. SELE and REG are the worst methods. The former passes from $\approx .252$ at $c = .99$ to $\approx .217$ at $c = .70$. The latter achieves $\approx .256$ at $c = .99$ and $c = .70$, with no improvement for small target coverages.

Next, we check the statistical significance of these results. For each target coverage and bootstrapped dataset, we rank the compared methods from 1 (the best) to 18 (the worst) w.r.t. the relative error rate. These rankings are then used in the Friedman’s omnibus test of equality of means and in its post-hoc Nemenyi test, following the steps described in Demsar (2006). Figure 5 shows Critical Difference (CD) plots, which provide a graphical representation of the output of the Nemenyi test. In each plot, the horizontal axis reports the average rank of each method – where being closer to one (farther to the right) implies better performances. A bold line connects methods whose differences are not statistically significant at 0.05 significance level. The plots show that there is no clear winner regardless of the coverage and of the binary/multiclass classification task. The group of not-statistically-different top methods contains between 8 and 14 baselines. However, ENS+SR is always the top ranked baseline, which makes it a good choice in general.

Q2. Comparing the empirical coverages. The constraint on the target coverage c in (9) is essential in many scenarios. Nevertheless, most papers do not sufficiently investigate the actual coverage achieved by the baselines. We assess how much the empirical coverage deviates from the target coverage c on the bootstrap dataset \mathcal{T}_{test} . To account for small coverage violations, we introduce a user-defined tolerance ε , and define the ε -coverage violation:

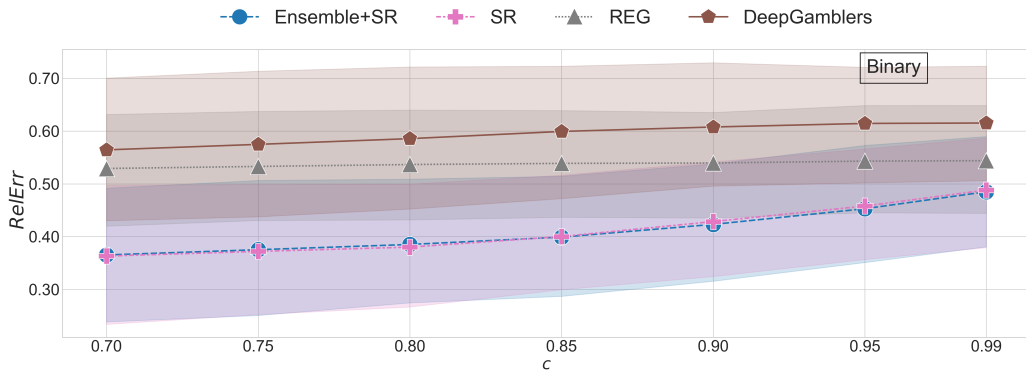
$$CovViol_{\varepsilon}(g, \mathcal{T}_{test}) = \min(0, \hat{\phi}(g, \mathcal{T}_{test}) - c + \varepsilon),$$

where $\hat{\phi}(g, \mathcal{T}_{test})$ is the empirical coverage on \mathcal{T}_{test} . Intuitively, $CovViol_{\varepsilon}$ is zero when the empirical coverage is greater or equal than the target coverage minus the tolerance; and it is greater than zero when the empirical coverage is smaller than $c - \varepsilon$. By looking at different tolerances ε , one can evaluate how the baselines perform w.r.t. small, medium or large coverage violations. We define the satisfaction of the constraint as:

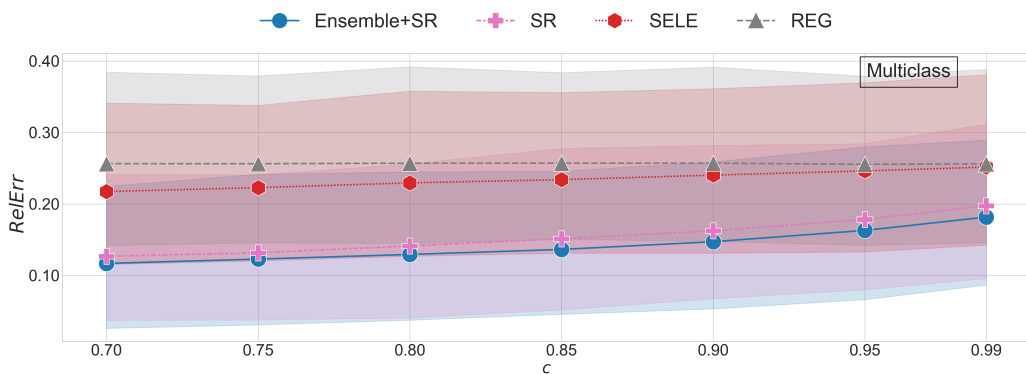
$$ConSat(\varepsilon) = \mathbb{1}(CovViol_{\varepsilon} = 0)$$

and report in Figure 6 the mean and standard deviation of $ConSat$ for ε in $\{0, .01, .02, .05, .10\}$. As for Figure 4, we limit the number of baselines to the two best and worst ones w.r.t. $ConSat$.

6. We rank baselines based on the mean value of relative error rate over all the target coverages.



(a) Binary datasets.



(b) Multiclass datasets.

Figure 4: Q1: $RelErr$ as a function of target coverage c for two best and worst approaches on (a) binary and (b) multiclass problems. On each subplot, only the two best and worst approaches are shown for readability.

As one would expect, the overall performances gradually improve for all baselines when increasing ε , and the gap among the baselines decreases. For binary data, the best methods are ENS and PLUGINAUC, and the worst methods are AUCROSS and SCROSS. When considering that no violation is allowed, i.e., $\varepsilon = 0$, the baselines satisfy the constraint between $\approx 39.9\%$ (CONFIDNET) and $\approx 56.5\%$ (SCROSS) of the times. For $\varepsilon = .01$ ENS has the highest value of $ConSat$ ($\approx .887$); for $\varepsilon = .02$ SAT is the best method (≈ 0.976); for $\varepsilon = .05$ both PLUGINAUC and SAT satisfy the constraint all the times.

For multiclass data, the top performers are SCROSS and SAT+EM, while the worst methods are REG and DG. At $\varepsilon = 0$, SCROSS has no coverage violations $\approx 75\%$ of the times, which is 25 percentage points more than the worst performing methods (SELE and REG). Interestingly, already at $\varepsilon = .02$, four methods (i.e., SCROSS, SAT+EM, SR, SAT+SR) always reach zero violations. For $\varepsilon = .05$, only CONFIDNET and SELNET+SR do not reach zero violations. In summary, coverage violations are generally limited, and noticeable differences among the baselines only occur at very small tolerances.

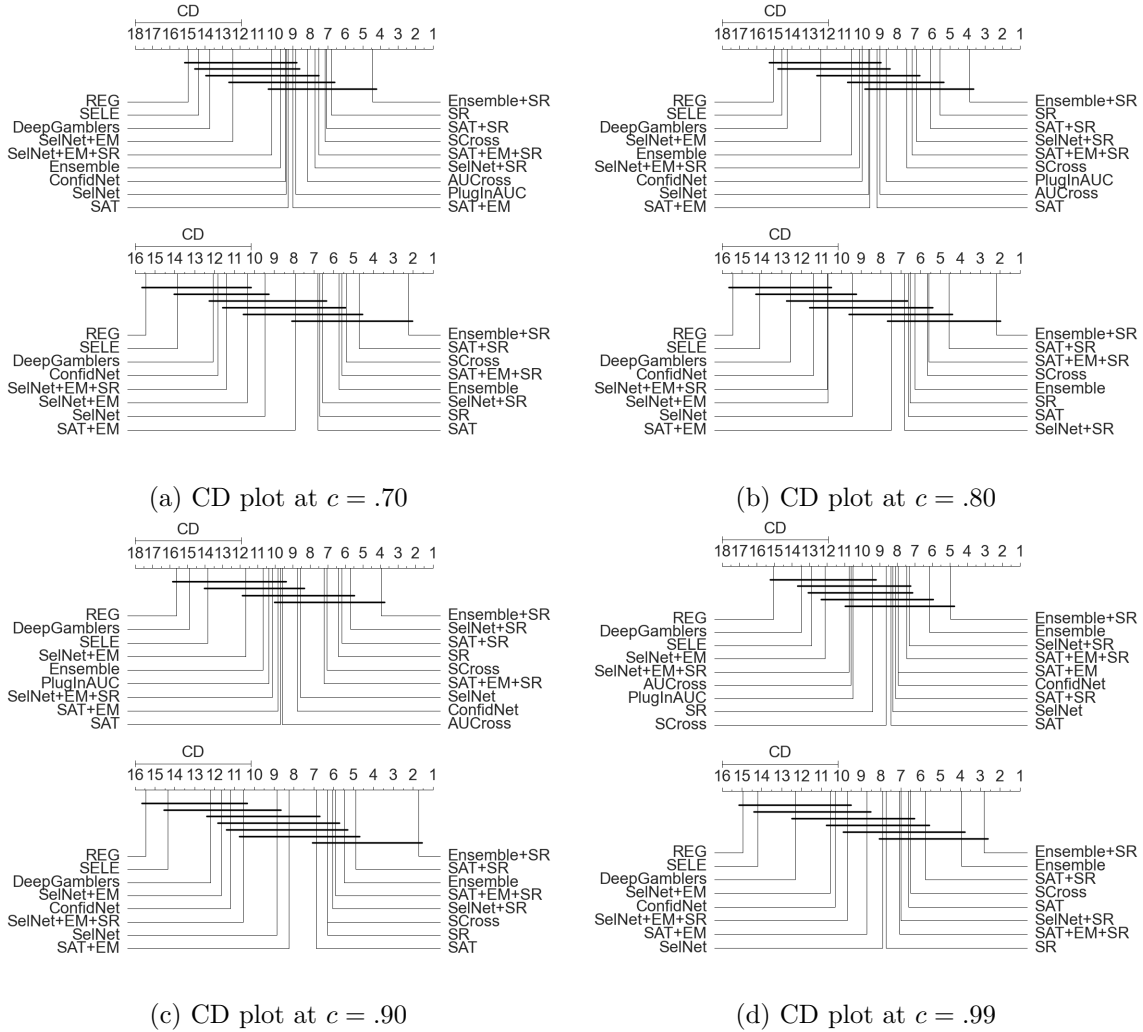
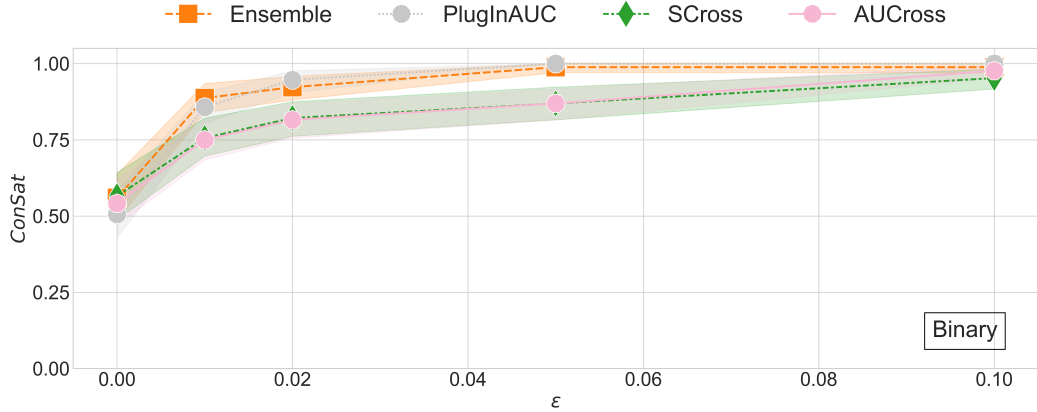


Figure 5: Q1: CD plots of relative error rate $RelErr$ for different target coverages. Top plots for binary datasets. Bottom plots for multiclass datasets.

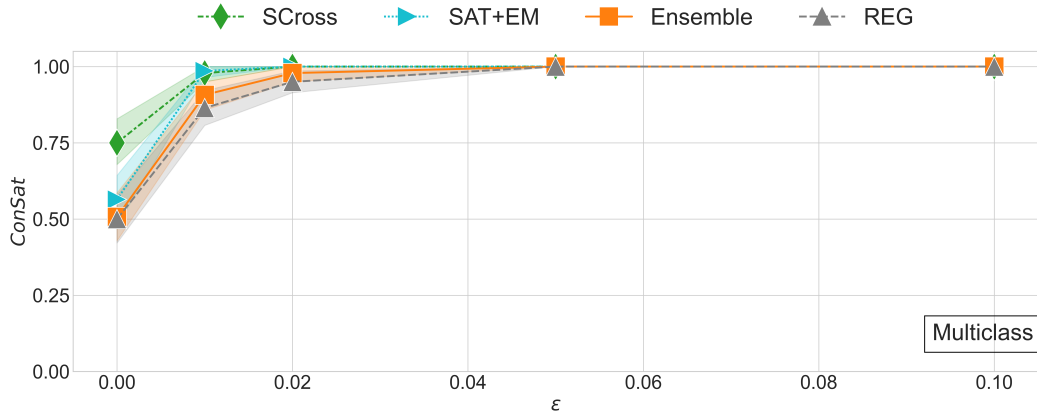
Q3. Rejection rate over classes. Pugnana and Ruggieri (2023b) observed that, in imbalanced classification tasks, selective classification methods reject proportionally more instances from the minority class. In this paragraph, we analyze this behavior on 7 binary class datasets of our collection with a minority class prior estimate $p \leq 0.25$ (Perini et al., 2020). Detailed results for the other binary datasets are reported in the Appendix B.3. First, let us introduce the *minority coefficient*:

$$MinCoeff = \frac{p_a}{p}, \tag{18}$$

defined as the ratio of the minority class proportion p_a in the accepted instances over the minority class prior p . Ideally, the minority coefficient should be ≈ 1 . Lower values indicate that the selective function introduces a bias against the minority class.



(a) Binary datasets.



(b) Multiclass datasets.

Figure 6: Q2: *ConSat* as a function of tolerance ϵ for two best and worst approaches on (a) binary and (b) multiclass problems. On each subplot, only the two best and worst approaches are shown for readability.

Figure 7 shows the mean minority coefficient for the best two and worst two baselines at the variation of the target coverage c . The best methods are AUCROSS and PLUGINAUC. Their minority coefficient is ≈ 1.00 and ≈ 1.01 respectively at $c = .99$, and it remains steady for lower coverages. At $c = .70$, PLUGINAUC reaches a mean *MinCoeff* ≈ 1.05 , and AUCROSS achieves *MinCoeff* ≈ 0.997 . For all the other baselines, there is a clear trend: the smaller the target coverage, the smaller the minority coefficient. For 11 out of 18 baselines, *MinCoeff* drops below .50 at $c = .70$. The worst methods are SELNET and DG. For the former, the mean *MinCoeff* ranges from $\approx .946$ at $c = .99$ to $\approx .375$ at $c = .70$. For the latter, the mean *MinCoeff* ranges from $\approx .966$ at $c = .99$ to $\approx .413$ for $c = .70$.

These results support the findings by Pugnana and Ruggieri (2023b), highlighting that the current approaches to SC, with the exception of AUCROSS and PLUGINAUC, do not take into account the issue of class balancing in the selected instances.

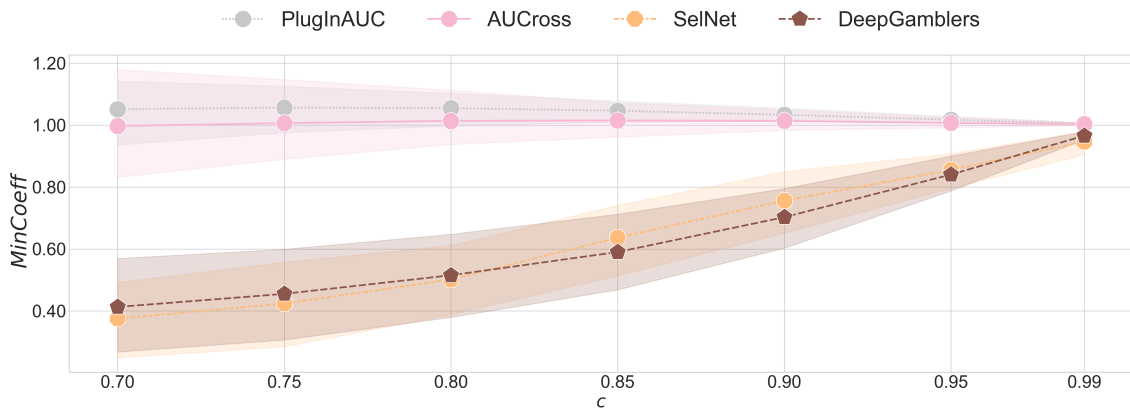


Figure 7: Q3: $MinCoeff$ as a function of target coverage c for two best and worst approaches. Only the two best and worst approaches are shown for readability.

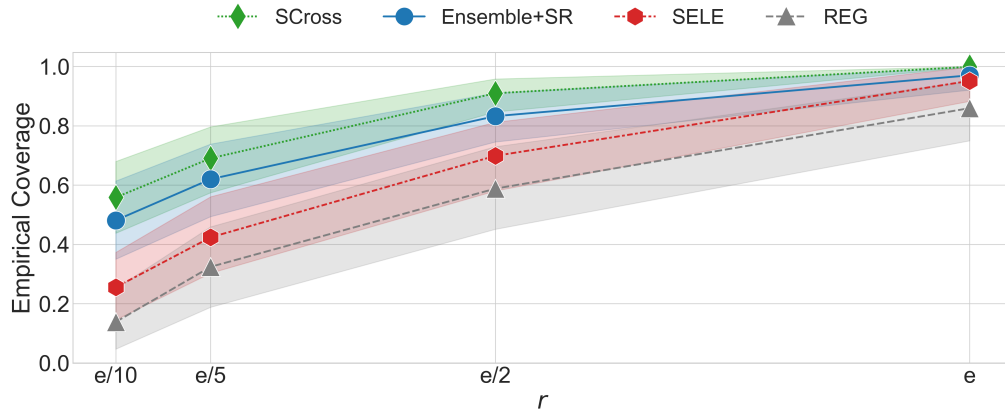
Q4. Flipping the learning task to maximize the model coverage under error constraints.

The vast majority of methods focus on the bounded-abstention model of Problem 2. To the best of our knowledge, the only method explicitly addressing the bounded-improvement model of Problem 1 is due to Gangrade et al. (2021), whose code has not been fully released. However, tackling the bounded-improvement model is useful in some application scenarios. Consider the bank example again. SC here can be used in two ways: on the one hand, the bank can set a target coverage c and calibrate a selective classifier so that $c\%$ of the cases are directly handled by the ML model, while the remaining - most difficult - ones are deferred to human experts. Here c is chosen on the basis of the personnel capacity of the bank. On the other hand, the bank can also be interested in maximizing the model coverage without incurring too many (costly) mistakes. Measuring such maximal coverage allows for planning the amount of human effort needed for the difficult cases.

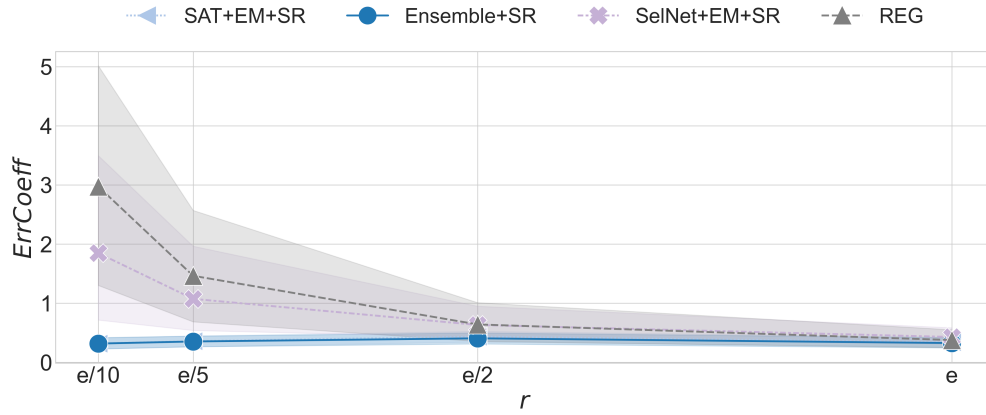
In this subsection, we evaluate the performances of the bounded-abstention baselines when flipping the task to the bounded-abstention problem through the *Selection with Guaranteed Risk* (SGR) algorithm proposed by Geifman and El-Yaniv (2017). SGR is a [classifier h , confidence k_h]-agnostic approach that optimizes the selection threshold τ (see (4)) such that the selective error rate at test time is guaranteed to be bounded ($\leq r$) with probability $> 1 - \delta$ and the coverage is maximized. We apply SGR on all the baselines but AUCROSS and PLUGINAUC, as their hard selection function is not compatible with SGR. Moreover, since SELNET needs specific coverage for training, we use all c 's one at a time, and compute the average results after applying SGR. We run experiments for four target error rates $r \in \{e/10, e/5, e/2, e\}$, where e is the dataset-specific error of the majority-class classifier h_{maj} on the whole test set, and set $\delta = 0.001$.

Figure 8 reports the results for the best two and the worst two baselines. The top plot shows the mean empirical coverage over the test sets of all baseline datasets (the higher, the better). The bottom plot shows the mean *error ratio* (the smaller, the better):

$$ErrCoeff = \hat{r}/r,$$



(a) W.r.t. the empirical coverage $\hat{\phi}$.



(b) W.r.t. the error ratio $ErrCoeff$.

Figure 8: Q4: SGR performance as a function of target error rate r for the two best and worst approaches in terms of (a) coverage $\hat{\phi}$, and (b) $ErrCoeff$. Only the two best and worst approaches are shown for readability.

between the empirical selective error rate \hat{r} and the target error rate r . When looking at the empirical coverage, the best performing baseline is SCROSS, with coverage ranging from .999 for $r = e$ to .558 for $r = e/10$. This is 40 percentage points higher than the worst method, namely REG. The second-best method is ENS+SR, with a mean coverage of $\approx .970$ at $r = e$ and $\approx .481$ at $r = e/10$.

Concerning $ErrCoeff$, we observe that for less strict target errors (i.e., e and $e/2$), all the baselines have error ratios close to 0. For more restrictive target errors, there is a gradual increase in the mean value of $ErrCoeff$. The methods with the smallest error ratios are ENS+SR and SAT+EM+SR, reaching $ErrCoeff \approx .316$ and $ErrCoeff \approx .317$ respectively at $r = e/10$. The worst methods are REG and SELNET+EM+SR, with a mean error ratio of ≈ 2.97 and ≈ 1.84 at $r = e/10$.

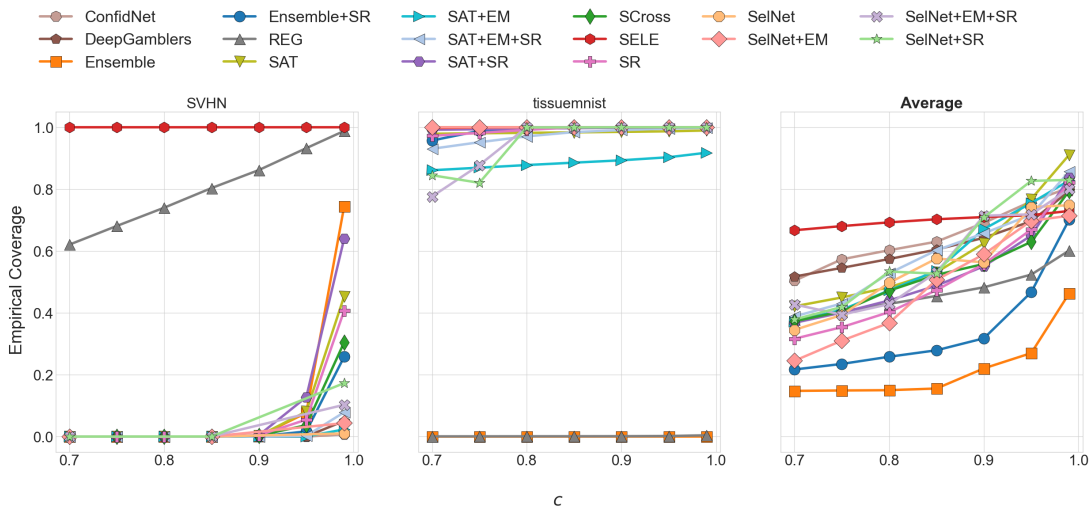


Figure 9: Q5: Empirical coverage $\hat{\phi}$ for out-of-distribution test sets (two selected image datasets and average results over the 20 image datasets) when varying coverages c .

Q5. Testing the methods on out-of-distribution examples. Although SC methods are not necessarily designed for working in out-of-distribution (o.o.d.) settings⁷, robustness of selective classifiers w.r.t. data shifts is highly sought. We investigate this property here by generating o.o.d. instances at test time for image datasets. We take an extreme approach by creating a test set with images made of uniformly random pixel values.⁸ An ideal selection function should reject the whole test set, since the images have close to zero probability of being drawn from the same distribution generating the training set.

For each of the 16 baselines that tackle multi-class classification, Figure 9 shows the empirical coverage obtained over the o.o.d. test set over the two selected image datasets and the average results over the 20 image datasets. Detailed results for all 20 datasets are provided in Appendix B.2.

First, we observe that there is no single method that manages to reject all the test instances across all datasets. For example, most of the baselines obtain empirical coverage close to 0 (best) for the 5 lowest coverage values on SVHN. On the other hand, on *tissuemnist*, the majority of baselines have always nearly maximum empirical coverage, with the sole exception of REG and ENS that manage to reject all the o.o.d. images.

Ensemble methods are generally better than other methods: ENS reaches the lowest mean empirical coverage on all datasets of $\approx .221$, ranging from $\approx .462$ at $c = .99$ to $\approx .148$ at $c = .70$, and ENS+SR is the second best method with a mean coverage $\approx .353$, ranging from $\approx .701$ at $c = .99$ to $\approx .217$ at $c = .70$.

5.3 Discussion

We conclude the experimental section by briefly summarizing the main findings.

7. See the novelty rejection approaches in the related work Section 6.

8. We provide additional results with less extreme shifts in the Appendix.

Regarding **Q1**, our results do not contradict the folk wisdom that an ensemble strategy (paired with the softmax response) overcome other baseline methods: ENS+SR always ranks first in terms of relative error rate. However, we stress that, depending on the target coverage, there are always at least nine baselines whose performance can’t be distinguished from ENS+SR’s one in a statistically significant sense. Conversely, some methods, i.e., REG, SELE, DG, SELNETSELNET+EM, SELNET+EM+SR, CONFIDNET, ENS, and PLUGINAUC, perform worse at least once (in a statistically significant sense) than top-performing baselines. Hence, our findings suggest that the claimed “superiority” of the considered state-of-the-art methods should be treated with caution: when we increase the number of experimental datasets, most methods perform equally well.

Our results on **Q2** show that coverage violations are generally small with a few exceptions of coverage violations above 10%. This confirms that the employed calibration strategies are well suited to achieving an empirical coverage that is fairly close to the target one.

For **Q3**, a significant difference arises among the methods. Only PLUGINAUC and AUCROSS reject equally across classes, while all the other methods abstain more relatively more frequently on the minority class. This behavior can have unforeseen consequences such as inducing cognitive bias in the human decision-maker that must make the decision on the rejected instances (Rastogi et al., 2022). For example, by abstaining more often on bad loan applicants, humans could be prone to associate the model’s rejections with bad applicants, even if this might not be necessarily true (Bondi et al., 2022).

The experiments for **Q4** show that SGR can effectively switch from the bounded abstention to the bounded improvement model assuming that target error rate is not too strict. In highly sensitive scenarios, where stronger guarantees are required, SGR often fails, thus suggesting a potential direction for future research towards methods specifically designed for the bounded improvement model.

For **Q5**, the results indicate that the current state-of-the-art baselines fail to reject consistently under distribution shifts. Consequently, practitioners should be cautious about applying SC techniques in the wild without considering potential issues deriving from data shifts. From a research perspective, this opens an intriguing future direction for shift-aware selective classification methods.

We point out that the methods which require training several neural networks might not be a feasible option for very large datasets, due to the huge computational power required. Such methods include the baselines ENS, ENS+SR, SCROSS, AUCROSS, CONFIDNET as well as the learn-to-select methods that require training a separate model for every target coverage.

6 Related Work

We present here a few related approaches and discuss in which respect they differ from SC.

Ambiguity Rejection. Ambiguity rejection focuses on abstaining on instances close to the decision boundary of the classifier (Hendrickx et al., 2024). SC is one of the main ways to perform ambiguity rejection. In particular, SC methods rely on confidence functions, which identify those instances where the classifier is more prone to make mistakes. Confidence values allow one to trade off coverage for selective risk.

The other main framework to perform ambiguity rejection is generally referred to as Learning to Reject (LtR) and is based on the seminal work by Chow (1970). Similarly to SC, LtR aims at learning a pair (classifier, rejector) such that the rejector determines when the classifier makes a prediction, limiting the predictions to the region where the classifier is likely correct (Cortes et al., 2023). However, LtR deviates from SC in two major aspects. First, the LtR methods learn the trade-off between abstention and prediction not by using confidence functions, but through a parameter a , representing the cost of rejection (Herbei and Wegkamp, 2006; Cortes et al., 2016; Tortorella, 2005; Condessa et al., 2013). However, setting the value of this hyperparameter is not straightforward, and it is context-dependant (Denis and Hebiri, 2020). Second, LtR methods are not meant to tackle the problem of minimizing a risk given a target coverage c . A more in-depth theoretical analysis for both LtR and SC can be found in (Franc et al., 2023), where the authors show that both frameworks share similar optimal strategies.

Novelty Rejection. A strategy orthogonal to ambiguity rejection consists of abstaining on instances that are unlikely to be seen according to the distribution of the training set. This approach is commonly referred to as *novelty rejection* (Dubuisson and Masson, 1993; Cordella et al., 1995), and is highly sought whenever there is a shift between the training and the test set distributions (Hendrickx et al., 2024; Van der Plas et al., 2023). Several techniques have been proposed for building novelty rejectors. As a first approach, one can estimate the marginal density and reject an instance if its probability is below a certain threshold (Nalisnick et al., 2019; Wang and Yiu, 2020). Another option is to employ a one-class classification model that predicts as novel the instances falling out of the region learnt from the training set (Coenen et al., 2020). Further approaches assign a score representing the novelty of an instance and abstain when such a score is above a certain level (Liang et al., 2018; Kühne et al., 2021; Perini and Davis, 2023; Van der Plas et al., 2023). To conclude, we highlight that the goal of novelty rejection differs from the SC goal, i.e. trading off risk and coverage, and linking the two problems is not straightforward (Hendrickx et al., 2024).

Conformal Prediction. Conformal prediction (Shafer and Vovk, 2008) augments the prediction of a M model by providing a set of target labels that comprise the true value with a specified (desired) level of confidence (Papadopoulos et al., 2002; Vovk, 2012; Kim et al., 2020; Abad et al., 2022; Angelopoulos et al., 2021). Differently from SC, conformal prediction focuses on quantifying the uncertainty associated with predictions rather than minimizing a specific type of error (Gangrade et al., 2021). Some works try to merge these two frameworks: for instance, in (Angelopoulos and Bates, 2021), conformal prediction is used to give guarantees over the selective error rate in an SC scenario by: (1) training a conformal predictor (e.g., SVC (Romano et al., 2020)), (2) calibrating its confidence levels, (3) setting a selection threshold over the confidence or p-values generated by the conformal predictor.

Learning to Defer. Learning to defer (Madras et al., 2018) is a generalization of LtR, where rather than incurring a rejection cost, the AI system can defer instances to human expert(s). One of the main differences in comparison to LtR and SC, is that the expert’s predictions might be wrong under the learning to defer framework. This is generally modelled using a cost function (Mozannar and Sontag, 2020). Thus, common methods include

the expert in the loop and aim to find an optimal assignment strategy for the whole human-AI system. Roughly speaking, such a strategy decides whether or not to make the model predict, which results in a cost equal to the model loss, or defer the prediction to the user, which incurs the user cost (Okati et al., 2021; De et al., 2020; Mozannar et al., 2023; Verma et al., 2023; Straitouri et al., 2022).

Real-world Applications. In recent years, abstaining AI systems have been deployed to foster human decision-making in increasingly many domains. For example, Van der Plas et al. (2023) describe a novelty rejector for sleep stage scoring. Cianci et al. (2023) exploit the SC strategy by Pugnana and Ruggieri (2023b) to augment a credit scoring ML model with an uncertainty self-assessment. Coenen et al. (2020) use unlabeled data on unaccepted loan applications to build a credit scoring model that can abstain from predicting. Hendrickx et al. (2021) propose a novelty rejector to find unexpected vehicle usage from sensor data and refrain from providing a prediction for such cases. Van Roy and Davis (2023) flag annotation errors in soccer data considering a specific confidence function for tree-based methods (Devos et al., 2023). Bondi et al. (2022) study a selective classifier deferring to humans to evaluate the presence of animals in photo traps. For other applications of abstaining classifiers, we refer to Hendrickx et al. (2024), while we refer to Punzi et al. (2024) for applications of hybrid-decision-making systems.

7 Conclusions

Limitations. For the sake of a fair comparison, our study focuses on neural network classifiers, as some of the methods assume a deep learning architecture for the classifier.

Due to the large computational costs of the experiments, for each dataset, we consider only a single deep-learning architecture chosen among the ones at the state-of-the-art. E.g., for `cifar10`, we implemented all the baselines using a VGG16 architecture. This might reduce the generalizability of our results to other deep-learning architectures.

We also acknowledge that a few studies, e.g., Gorishniy et al. (2021); Grinsztajn et al. (2022), point out that for tabular datasets, the usage of tree-based models is the current state of the art. In this sense, model-agnostic methods could benefit from using other base classifiers, as shown in Pugnana and Ruggieri (2023a,b).

Another limitation of our benchmark is that we consider only images and tabular data, since they are the main data type over which SC methods have been tested so far. This choice is in line with the goal of this paper, which aims to compare existing approaches fairly. However, our results do not necessarily extend to other kinds of data such as text, audio or time series.

A possible concern could also regard the size of the datasets in our benchmark, which never exceeds $\approx 300k$ instances. This aspect could impact the external validity of our discussion. However, there are reasons for this choice. First, the considered datasets are used either in popular benchmarks (Yang et al., 2023; Gorishniy et al., 2021; Grinsztajn et al., 2022), or by selective classification works, e.g., (Geifman and El-Yaniv, 2019; Franc et al., 2023; Pugnana and Ruggieri, 2023b). Second, since we trained and fine-tuned all the models from scratch, with considerable computational costs, we decided to prioritize the variety of data over the size of a single dataset.

Moreover, our bootstrap procedure quantifies variability only in the test set. According to several works, such as Kohavi (1995), the best resampling method is stratified k-fold cross-validation with $K = 10$. Unfortunately, these strategies are not computationally feasible when employing large neural networks as in our study. Hence, we had to opt for a single train-test-split and bootstrap only the test set, as done for instance by Rajkomar et al. (2018).

Finally, our study does not report on the running times of the baselines, since, due to load balancing issues, we had to distribute the experiments over several machines with different hardware settings. This made it impossible to compare the running times of runs over different machines. However, we point out that the running times are proportional to the number of training tasks required by each method. E.g., ENS requires to train ten neural networks (see Appendix A.2), leading to a running time of about ten times the one of PLUGINAUC, which requires to train a single neural network.

Conclusions. We extensively evaluated 18 SC baselines over 44 datasets, taking into account both images and tabular data as well as both binary and multiclass classification tasks. Regarding previously investigated tasks, our extended analysis shows that: *(i)* there are no statistically significant differences among most of the methods in terms of selective error rate, even though ENS+SR always ranks first across all the baselines; *(ii)* large coverage violations are rare for all the methods with no significant difference among baselines for our data; *(iii)* on binary classification tasks, we observed different patterns between imbalanced and balanced domains regarding rejection rates across classes: in the former case, only AUCROSS and PLUGINAUC succeeded in not primarily rejecting the minority class. Moreover, we also emphasize novel findings: *(iv)* we tested empirically the effectiveness of SGR to switch from the bounded-abstention setting to the bounded-improvement one, noticing room for improvement when a very small target error rate is required; *(v)* we show how current methods fail in correctly rejecting instances when extreme feature shifts occur, pointing to a highly relevant open problem in the area.

Broader Impact Statement

Because Machine Learning models can make errors in their predictions, adding a reject option is a means for improving their trustworthiness. The selective classification framework is one of the most popular ways to achieve such a goal by coupling a classifier with a selective function that decides whether to accept or reject making a prediction. However, existing selective classification methods have never been evaluated on a large scale. Our work is the first to fill this gap, providing the first extensive benchmark for testing selective classification methods. The experimental evaluation sheds light on the strengths and weaknesses of selective classification methods for what concerns their error rate, acceptance rate (called coverage), distribution of rejection over classes, and robustness to data shifts.

Acknowledgments and Disclosure of Funding

The work of A. Pugnana and S. Ruggieri has been partly funded by PNRR - M4C2 - Investimento 1.3, Partenariato Esteso PE00000013 - "FAIR - Future Artificial Intelligence

Research” - Spoke 1 “Human-centered AI”, funded by the European Commission under the NextGeneration EU programme, and by the project FINDHR funded by the EU’s Horizon Europe research and innovation program under g.a. No 101070212. Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the EU. Neither the EU nor the granting authority can be held responsible for them.

L. Perini received funding from FWOVlaanderen (aspirant grant 1166222N). Moreover, L. Perini and J. Davis received funding from the Flemish Government under the “Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen” programme.

References

- Javier Abad, Umang Bhatt, Adrian Weller, and Giovanni Cherubin. Approximating full conformal prediction at scale via influence functions. *CoRR*, abs/2202.01315, 2022.
- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *KDD*, pages 2623–2631. ACM, 2019.
- Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *CoRR*, abs/2107.07511, 2021.
- Anastasios Nikolas Angelopoulos, Stephen Bates, Michael I. Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. In *ICLR*. OpenReview.net, 2021.
- Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. Botorch: A framework for efficient monte-carlo bayesian optimization. In *NeurIPS*, 2020.
- Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: The state of the art. *Sociological Methods & Research*, 50(1):3–44, 2021.
- Elizabeth Bondi, Raphael Koster, Hannah Sheahan, Martin J. Chadwick, Yoram Bachrach, A. Taylan Cemgil, Ulrich Paquet, and Krishnamurthy Dvijotham. Role of human-AI interaction in selective prediction. In *AAAI*, pages 5286–5294. AAAI Press, 2022.
- C. K. Chow. On optimum recognition error and reject tradeoff. *IEEE Trans. Inf. Theory*, 16(1):41–46, 1970.
- Giuseppe Cianci, Roberto Goglia, Riccardo Guidotti, Matteo Kapllaj, Roberto Mosca, Andrea Pugnana, Franco Ricotti, and Salvatore Ruggieri. Applied data science for leasing score prediction. In *IEEE Big Data*, pages 1687–1696. IEEE, 2023.
- Lize Coenen, Ahmed K. A. Abdullah, and Tias Guns. Probability of default estimation, with a reject option. In *DSAA*, pages 439–448. IEEE, 2020.
- Filipe Condessa, José M. Bioucas-Dias, Carlos A. Castro, John A. Ozolek, and Jelena Kovacevic. Classification with reject option using contextual information. In *ISBI*, pages 1340–1343. IEEE, 2013.

- Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. Addressing failure prediction by learning model confidence. In *NeurIPS*, pages 2898–2909, 2019.
- Luigi P. Cordella, Claudio De Stefano, Carlo Sansone, and Mario Vento. An adaptive reject option for LVQ classifiers. In *ICIAP*, volume 974 of *Lecture Notes in Computer Science*, pages 68–73. Springer, 1995.
- Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Boosting with abstention. In *NeurIPS*, pages 1660–1668, 2016.
- Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Theory and algorithms for learning with rejection in binary classification. *Annals of Mathematics and Artificial Intelligence*, pages 1–39, 2023.
- Sarah JC Craig, Ana M Kenney, Junli Lin, Ian M Paul, Leann L Birch, Jennifer S Savage, Michele E Marini, Francesca Chiaromonte, Matthew L Reimherr, and Kateryna D Makova. Constructing a polygenic risk score for childhood obesity using functional data analysis. *Econometrics and Statistics*, 25:66–86, 2023.
- Xolani Dastile, Turgay Çelik, and Moshe Potsane. Statistical and machine learning models in credit scoring: A systematic literature survey. *Appl. Soft Comput.*, 91:106263, 2020.
- Abir De, Paramita Koley, Niloy Ganguly, and Manuel Gomez-Rodriguez. Regression under human assistance. In *AAAI*, pages 2611–2620. AAAI Press, 2020.
- Janez Demsar. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn. Res.*, 7:1–30, 2006.
- Christophe Denis and Mohamed Hebiri. Consistency of plug-in confidence sets for classification in semi-supervised learning. *J. of Nonpar. Statistics*, 32(1):42–72, 2020.
- Laurens Devos, Lorenzo Perini, Wannes Meert, and Jesse Davis. Detecting evasion attacks in deployed tree ensembles. In *ECML/PKDD (5)*, volume 14173 of *Lecture Notes in Computer Science*, pages 120–136. Springer, 2023.
- Qiang Ding, Yixuan Cao, and Ping Luo. Top-ambiguity samples matter: Understanding why deep ensemble works in selective classification. In *NeurIPS*, 2023.
- Bernard Dubuisson and Mylène Masson. A statistical decision rule with incomplete knowledge about classes. *Pattern Recognit.*, 26(1):155–165, 1993.
- Ran El-Yaniv and Yair Wiener. On the foundations of noise-free selective classification. *J. Mach. Learn. Res.*, 11:1605–1641, 2010.
- European Commission. Proposal for a Regulation of the European Parliament and of the Council Laying down harmonised rules on Artificial Intelligence (AI Act) and amending certain Union legislative acts, 2021. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206>.

- Alessandro Fabris, Nina Baranowska, Matthew J. Dennis, Philipp Hacker, Jorge Saldivar, Frederik J. Zuiderveen Borgesius, and Asia J. Biega. Fairness and bias in algorithmic hiring. *CoRR*, abs/2309.13933, 2023.
- Leo Feng, Mohamed Osama Ahmed, Hossein Hajimirsadeghi, and Amir H. Abdi. Towards better selective classification. In *ICLR*. OpenReview.net, 2023.
- Vojtech Franc, Daniel Průša, and Václav Voráček. Optimal strategies for reject option classifiers. *J. Mach. Learn. Res.*, 24:11:1–11:49, 2023.
- Aditya Gangrade, Anil Kag, and Venkatesh Saligrama. Selective classification via one-sided prediction. In *AISTATS*, volume 130, pages 2179–2187. PMLR, 2021.
- Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In *NIPS*, pages 4878–4887, 2017.
- Yonatan Geifman and Ran El-Yaniv. Selectivenet: A deep neural network with an integrated reject option. In *ICML*, volume 97, pages 2151–2159. PMLR, 2019.
- Yury Gorishniy, Ivan Rubachev, Valentin Khruikov, and Artem Babenko. Revisiting deep learning models for tabular data. In *NeurIPS*, pages 18932–18943, 2021.
- Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? In *NeurIPS*, 2022.
- Haibo He and Edwardo A. Garcia. Learning from imbalanced data. *IEEE Trans. Knowl. Data Eng.*, 21(9):1263–1284, 2009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778. IEEE Computer Society, 2016.
- Kilian Hendrickx, Wannes Meert, Bram Cornelis, and Jesse Davis. Know your limits: Machine learning with rejection for vehicle engineering. In *ADMA*, volume 13087 of *Lecture Notes in Computer Science*, pages 273–288. Springer, 2021.
- Kilian Hendrickx, Lorenzo Perini, Dries Van der Plas, Wannes Meert, and Jesse Davis. Machine learning with a reject option: a survey. *Mach. Learn.*, 113(5):3073–3110, 2024.
- Radu Herbei and Maten H. Wegkamp. Classification with reject option. *Can. J. Stat.*, 34(4):709–721, 2006.
- Lang Huang, Chao Zhang, and Hongyang Zhang. Self-adaptive training: beyond empirical risk minimization. In *NeurIPS*, 2020.
- Erik Jones, Shiori Sagawa, Pang Wei Koh, Ananya Kumar, and Percy Liang. Selective classification can magnify disparities across groups. In *ICLR*. OpenReview.net, 2021.
- Davinder Kaur, Suleyman Uslu, Kaley J. Rittichier, and Arjan Durrresi. Trustworthy Artificial Intelligence: A review. *ACM Comput. Surv.*, 55(2):39:1–39:38, 2023.

- Byol Kim, Chen Xu, and Rina Foygel Barber. Predictive inference is free with the jackknife+-after-bootstrap. In *NeurIPS*, 2020.
- Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *IJCAI*, pages 1137–1145. Morgan Kaufmann, 1995.
- Joana Kühne, Christian März, et al. Securing deep learning models with autoencoder based anomaly detection. In *PHM Society European Conference*, volume 6, pages 13–13, 2021.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NIPS*, pages 6402–6413, 2017.
- Matilde Lazzari, José M. Álvarez, and Salvatore Ruggieri. Predicting and explaining employee turnover intention. *Int. J. Data Sci. Anal.*, 14(3):279–292, 2022.
- Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *ICLR*. OpenReview.net, 2018.
- Ziyin Liu, Zhikang Wang, Paul Pu Liang, Ruslan Salakhutdinov, Louis-Philippe Morency, and Masahito Ueda. Deep gamblers: Learning to abstain with portfolio theory. In *NeurIPS*, pages 10622–10632, 2019.
- David Madras, Toniann Pitassi, and Richard S. Zemel. Predict responsibly: Improving fairness and accuracy by learning to defer. In *NeurIPS*, pages 6150–6160, 2018.
- Hussein Mozannar and David A. Sontag. Consistent estimators for learning to defer to an expert. In *ICML*, volume 119, pages 7076–7087. PMLR, 2020.
- Hussein Mozannar, Hunter Lang, Dennis Wei, Prasanna Sattigeri, Subhro Das, and David A. Sontag. Who should predict? exact algorithms for learning to defer to humans. In *AISTATS*, volume 206, pages 10520–10545. PMLR, 2023.
- Eric T. Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Görür, and Balaji Lakshminarayanan. Hybrid models with deep and invertible features. In *ICML*, volume 97, pages 4723–4732. PMLR, 2019.
- Nastaran Okati, Abir De, and Manuel Gomez-Rodriguez. Differentiable learning under triage. In *NeurIPS*, pages 9140–9151, 2021.
- Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alexander Gammerman. Inductive confidence machines for regression. In *ECML*, volume 2430 of *Lecture Notes in Computer Science*, pages 345–356. Springer, 2002.
- Lorenzo Perini and Jesse Davis. Unsupervised anomaly detection with rejection. In *NeurIPS*, 2023.
- Lorenzo Perini, Vincent Vercruyssen, and Jesse Davis. Class prior estimation in active positive and unlabeled learning. In *IJCAI*, pages 2915–2921. ijcai.org, 2020.
- Andrea Pugnana. Topics in selective classification. In *AAAI*, pages 16129–16130. AAAI Press, 2023.

- Andrea Pugnana and Salvatore Ruggieri. A model-agnostic heuristics for selective classification. In *AAAI*, pages 9461–9469. AAAI Press, 2023a.
- Andrea Pugnana and Salvatore Ruggieri. AUC-based selective classification. In *AISTATS*, volume 206, pages 2494–2514. PMLR, 2023b.
- Clara Punzi, Roberto Pellungrini, Mattia Setzu, Fosca Giannotti, and Dino Pedreschi. Ai, meet human: Learning paradigms for hybrid decision making systems. *CoRR*, abs/2402.06287, 2024.
- Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Michaela Hardt, Peter J Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, et al. Scalable and accurate deep learning with electronic health records. *NPJ Digital Medicine*, 1(1):1–10, 2018.
- Charvi Rastogi, Yunfeng Zhang, Dennis Wei, Kush R. Varshney, Amit Dhurandhar, and Richard Tomsett. Deciding fast and slow: The role of cognitive biases in AI-assisted decision-making. *Proc. ACM Hum. Comput. Interact.*, 6(CSCW1):83:1–83:22, 2022.
- Yaniv Romano, Matteo Sesia, and Emmanuel J. Candès. Classification with valid and adaptive coverage. In *NeurIPS*, 2020.
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *J. Mach. Learn. Res.*, 9:371–421, 2008.
- Eleni Straitouri, Lequn Wang, Nastaran Okati, and Manuel Gomez Rodriguez. Provably improving expert predictions with conformal prediction. *CoRR*, abs/2201.12006, 2022.
- Francesco Tortorella. A ROC-based reject rule for dichotomizers. *Pattern Recognit. Lett.*, 26(2):167–180, 2005.
- Dries Van der Plas, Wannes Meert, Johan Verbraecken, and Jesse Davis. A novel reject option applied to sleep stage scoring. In *SDM*, pages 820–828. SIAM, 2023.
- Maaïke Van Roy and Jesse Davis. Datadebugging: Enhancing trust in soccer action-value models by contextualization. In *13th World Congress of Performance Analysis of Sport and 13th International Symposium on Computer Science in Sport*, pages 193–196, 2023. ISBN 978-3-031-31772-9.
- Rajeev Verma, Daniel Barrejón, and Eric T. Nalisnick. Learning to defer to multiple experts: Consistent surrogate losses, confidence calibration, and conformal ensembles. In *AISTATS*, volume 206, pages 11415–11434. PMLR, 2023.
- Vladimir Vovk. Cross-conformal predictors. *CoRR*, abs/1208.0806, 2012.
- Xin Wang and Siu-Ming Yiu. Classification with rejection: Scaling generative classifiers with supervised deep infomax. In *IJCAI*, pages 2980–2986. ijcai.org, 2020.
- Jiancheng Yang, Rui Shi, Donglai Wei, Zequan Liu, Lin Zhao, Bilian Ke, Hanspeter Pfister, and Bingbing Ni. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1):41, 2023.

Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyu Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan Hendrycks, Yixuan Li, and Ziwei Liu. Openood: Benchmarking generalized out-of-distribution detection. In *NeurIPS*, 2022.

Tianbao Yang and Yiming Ying. AUC maximization in the era of big data and AI: A survey. *ACM Comput. Surv.*, 55(8):172:1–172:37, 2023.

Appendix A. Experimental Details

We provide all the additional information on datasets, settings, and code for replicating the experiments of the paper.

A.1 Datasets

Table A1 reports the datasets used in our benchmarking and a link to retrieve the original data. We also include whether the dataset was considered in a previous SC evaluation.

Table A2 reports some experimental details, including training size; batch size used for training; feature space in terms of features for tabular data and image size for image data; target space number of classes m ; the percentage of the minority class in each dataset. We also report the Deep Neural Network (DNN) architectures we employed for each dataset. Such a choice was made according to the following criteria:

- (i) a former paper in the literature used this dataset and employed a specific architecture;
- (ii) if point (i) does not apply, we did the following:
 - for image data, we employed a ResNet34 architecture;
 - for tabular data, if a dataset was tested in Gorishniy et al. (2021), we applied the best-performing architecture on that specific dataset. Otherwise, we employed the FTTransformer architecture following the suggestion by Grinsztajn et al. (2022).

All the data were re-shuffled, normalized and split into training, test, calibration and validation sets, according to a 60%, 20%, 10%, and 10% proportion (respectively). In the code repository, we provide the Python scripts to recreate the data employed in this analysis.

A.2 Hyperparameter Settings

We optimize the hyperparameters using `Optuna` (Akiba et al., 2019), a framework for multi-objective Bayesian optimization, with the following inputs: coverage violation and cross-entropy loss as target metrics, `BoTorch` as sampler (Balandat et al., 2020), 10 initial independent trials out of 20 total trials. We report in Table A3 the parameter space we used during the tuning procedure. Some hyperparameters are loss-specific, as they refer to a specific baseline loss, while others are network-specific, as they refer to a specific deep neural network architecture. We report the search space in the last column, where the notation $[a; b]_X$ stands for all values within $[a; b]$ that are linearly spaced with a gap equal to X . For example, $[0, 60]_{15}$ indicates the set of values $\{0, 15, 30, 45, 60\}$. For small sets of values, we directly indicate all the possibilities using the notation $\{a_1, a_2, \dots\}$. We discretize the search space of most hyperparameters because Bayesian Optimization suffers from high-dimensional spaces. We specify that if `time decay` is set to `True`, we halve the learning rate every 25 epochs as done in Geifman and El-Yaniv (2019). Moreover, for image data, which we found to be more unstable during the training, we start the `Optuna` optimization procedure using default values if these were suggested in some SC paper. Due to the huge computational cost of the tuning procedure, `SELNET+SR`, `SELNET+EM+SR`, `SAT+SR` and `SAT+EM+SR` use the same optimal hyperparameters as,

Table A1: Dataset sources.

| Dataset | Data Type | Link | Previous SC paper |
|----------------|-----------|-------------------------------------|--|
| adult | Tabular | uci/adult | Pugnana and Ruggieri (2023a,b) |
| aloi | Tabular | openml/id=1592 | – |
| bank | Tabular | uci/bank+marketing | Franc et al. (2023) |
| bloodmnist | Image | zenodo/6496656/files/bloodmnist | – |
| breastmnist | Image | zenodo/6496656/files/breastmnist | – |
| catsdogs | Image | kaggle/dogs-vs-cats | Geifman and El-Yaniv (2019); Liu et al. (2019); Huang et al. (2020); Pugnana and Ruggieri (2023a,b) |
| chestmnist | Image | zenodo/6496656/files/chestmnist | – |
| cifar10 | Image | pytorch/vision/CIFAR10 | Geifman and El-Yaniv (2019); Corbière et al. (2019); Huang et al. (2020); Feng et al. (2023); Pugnana and Ruggieri (2023b) |
| compass | Tabular | openml/id=44162 | – |
| covtype | Tabular | openml/id=1596 | Franc et al. (2023) |
| dermannist | Image | zenodo/6496656/files/dermannist | – |
| electricity | Tabular | openml/id=44120 | – |
| eye | Tabular | openml/id=44157 | – |
| food101 | Image | pytorch/vision/Food101 | Feng et al. (2023) |
| giveme | Tabular | kaggle/GiveMeSomeCredit | Pugnana and Ruggieri (2023a,b) |
| helena | Tabular | openml/id=41169 | – |
| heloc | Tabular | openml/id=45023 | – |
| higgs | Tabular | openml/id=23512 | – |
| house | Tabular | openml/id=43957 | – |
| indian | Tabular | openml/id=41972 | – |
| jannis | Tabular | openml/id=44079 | – |
| kddipums97 | Tabular | openml/id=44124 | – |
| letter | Tabular | openml/id=6 | Franc et al. (2023) |
| magic | Tabular | openml/id=44125 | – |
| miniboone | Tabular | openml/id=44119 | – |
| MNIST | Image | pytorch/vision/MNIST | Lakshminarayanan et al. (2017); Liu et al. (2019); Corbière et al. (2019) |
| octmnist | Image | zenodo/6496656/files/octmnist | – |
| online | Tabular | openml/id=45060 | – |
| organamnist | Image | zenodo/6496656/files/organamnist | – |
| organcmnist | Image | zenodo/6496656/files/organcmnist | – |
| organsmnist | Image | zenodo/6496656/files/organsmnist | – |
| oxfordpets | Image | pytorch/vision/oxfordpets | – |
| pathmnist | Image | zenodo/6496656/files/pathmnist | – |
| phoneme | Tabular | openml/id=44127 | – |
| pneumoniamnist | Image | zenodo/6496656/files/pneumoniamnist | – |
| pol | Tabular | openml/id=43991 | – |
| retinamnist | Image | zenodo/6496656/files/retinamnist | – |
| rl | Tabular | openml/id=44160 | – |
| stanfordcars | Image | pytorch/vision/StanfordCars | Feng et al. (2023) |
| SVHN | Image | pytorch/vision/SVHN | Geifman and El-Yaniv (2017, 2019); Liu et al. (2019); Corbière et al. (2019) |
| tissuemnist | Image | zenodo/6496656/files/tissuemnist | – |
| ucicredit | Tabular | openml/id=42477 | Pugnana and Ruggieri (2023a,b) |
| upselling | Tabular | openml/id=44158 | – |
| waterbirds | Image | stanford.edu/dro/waterbird | Jones et al. (2021) |

Table A2: Dataset details.

| Dataset | Training Size | Batch Size | # Features | # Classes | Minority Ratio | DNN Architecture |
|----------------|---------------|------------|------------------|-----------|----------------|------------------|
| adult | 29,303 | 256 | 13 | 2 | 23.9% | FTTransformer |
| aloi | 64,800 | 512 | 128 | 1000 | 0.1% | TabResnet |
| bank | 27,126 | 128 | 16 | 2 | 11.7% | TabResnet |
| bloodmnist | 10,253 | 128 | 28×28 | 8 | 7.1% | Resnet18 |
| breastmnist | 468 | 64 | 28×28 | 2 | 26.9% | Resnet18 |
| catsdogs | 15,000 | 128 | 64×64 | 2 | 50.0% | VGG |
| chestmnist | 67,272 | 512 | 28×28 | 2 | 10.3% | Resnet18 |
| cifar10 | 36,000 | 128 | 32×32 | 1 | 100.0% | VGG |
| compass | 9,985 | 128 | 17 | 2 | 50.0% | FTTransformer |
| covtype | 348,605 | 1024 | 54 | 7 | 0.5% | FTTransformer |
| dermamnist | 6,008 | 128 | 28×28 | 7 | 1.1% | Resnet18 |
| electricity | 23,083 | 128 | 7 | 2 | 50.0% | FTTransformer |
| eye | 4,564 | 128 | 23 | 2 | 50.0% | FTTransformer |
| food101 | 60,600 | 256 | 224×224 | 101 | 1.0% | Resnet34 |
| giveme | 90,000 | 512 | 8 | 2 | 6.7% | TabResnet |
| helenas | 39,116 | 512 | 27 | 100 | 0.2% | TabResnet |
| heloc | 6,000 | 128 | 22 | 2 | 50.0% | FTTransformer |
| higgs | 58,829 | 512 | 28 | 2 | 47.1% | FTTransformer |
| house | 8,092 | 128 | 16 | 2 | 50.0% | FTTransformer |
| indian | 5,485 | 128 | 220 | 8 | 0.2% | TabResnet |
| jannis | 34,548 | 512 | 54 | 2 | 50.0% | FTTransformer |
| kddipums97 | 3,112 | 128 | 20 | 2 | 50.0% | FTTransformer |
| letter | 12,000 | 128 | 16 | 26 | 3.7% | FTTransformer |
| magic | 8,024 | 128 | 10 | 2 | 50.0% | FTTransformer |
| miniboone | 43,798 | 256 | 50 | 2 | 50.0% | TabResnet |
| MNIST | 42,000 | 128 | 28×28 | 10 | 9.0% | Resnet34 |
| octmnist | 65,585 | 512 | 28×28 | 4 | 8.1% | Resnet18 |
| online | 7,398 | 128 | 17 | 2 | 15.5% | FTTransformer |
| organamnist | 35,310 | 256 | 28×28 | 11 | 4.0% | Resnet18 |
| organcmnist | 14,196 | 128 | 28×28 | 11 | 4.8% | Resnet18 |
| organsmnist | 15,132 | 128 | 28×28 | 11 | 4.6% | Resnet18 |
| oxfordpets | 4,409 | 128 | 224×224 | 2 | 32.3% | Resnet34 |
| pathmnist | 64,308 | 512 | 28×28 | 9 | 8.9% | Resnet18 |
| phoneme | 1,901 | 128 | 5 | 2 | 50.0% | FTTransformer |
| pneumoniamnist | 3,512 | 128 | 28×28 | 2 | 27.1% | Resnet18 |
| pol | 9,000 | 128 | 26 | 11 | 1.7% | FTTransformer |
| retinamnist | 960 | 128 | 28×28 | 5 | 5.8% | Resnet18 |
| rl | 2,982 | 128 | 12 | 2 | 50.0% | FTTransformer |
| stanfordcars | 9,710 | 128 | 224×224 | 196 | 0.3% | Resnet34 |
| SVHN | 59,573 | 128 | 32×32 | 10 | 6.3% | VGG |
| tissuemnist | 141,830 | 1024 | 28×28 | 8 | 3.5% | Resnet18 |
| ucicredit | 18,000 | 128 | 23 | 2 | 22.1% | TabResnet |
| upselling | 3,017 | 128 | 45 | 2 | 50.0% | FTTransformer |
| waterbirds | 7,072 | 128 | 224×224 | 2 | 22.6% | Resnet50 |

respectively, SELNET, SELNET+EM, SAT and SAT+EM. Similarly, SCROSS, ENS, ENS+SR, AUCROSS, and PLUGINAUC employ the best configuration found for SR as they share the same training loss, i.e., cross-entropy. For both SCROSS and AUCROSS we set $K = 5$, following the suggestions in Pugnana and Ruggieri (2023a,b). For both ENS and ENS+SR we used the default value of $K = 10$, following the suggestions in Lakshminarayanan et al. (2017). For the uncertainty network of CONFIDNET, we employed the same choice architecture detailed in the original paper (Corbière et al., 2019), i.e., the same main body as the network classifier followed by 4 dense layers in a single node with sigmoid activation. We used such a structure also for building SELE and REG uncertainty

Table A3: Hyperparameter spaces.

| Parameter | Loss-Specific | Network-Specific | Search Space |
|--------------------|----------------------|--------------------------|--|
| o | DG | No | $[1; m]_{.02}$ |
| γ | SAT, SAT+EM | No | $[0.9; 0.99]_{.01}$ |
| E_s | SAT, SAT+EM | No | $[0; 60]_{15}$ |
| β | SAT+EM, SELNET+EM | No | $\{1e-4, 1e-3, 1e-2, 1e-1\}$ |
| α | SELNET, SELNET+EM | No | $\{.25; .75\}_{.05}$ |
| λ | SELNET, SELNET+EM | No | $\{8, 16, 32, 64\}$ |
| optimizer | No | No | $\{\text{SGD, Adam, AdamW}\}$ |
| learning rate | No | No | $\{1e-5, 1e-4, 1e-3, 1e-2\}$ |
| optimizer unc. | CONFIDNET, SELE, REG | Yes | $\{\text{SGD, Adam, AdamW}\}$ |
| learning rate unc. | CONFIDNET, SELE, REG | Yes | $\{1e-8, 1e-7, 1e-6, 1e-5, 1e-4, 1e-3, 1e-2\}$ |
| time decay | No | No | $\{\text{True, False}\}$ |
| nesterov | No | No | $\{\text{True, False}\}$ |
| nesterov unc. | CONFIDNET, SELE, REG | Yes | $\{\text{True, False}\}$ |
| weight decay | No | No | $\{1e-6, 1e-5, 1e-4, 1e-3\}$ |
| d_token | No | FTTransformer, TabResNet | $[64, 512]_{64}$ |
| n_blocks | No | FTTransformer, TabResNet | $\{1, 2, 3, 4\}$ |
| d_hidden_factor | No | FTTransformer, TabResNet | $[2/3; 8/3]_{1/3}$ |
| attention_dropout | No | FTTransformer | $\{0; .5\}_{.05}$ |
| residual_dropout | No | FTTransformer | $\{0; .2\}_{.05}$ |
| ffn_dropout | No | FTTransformer | $\{0; .5\}_{.05}$ |
| d_main | No | TabResNet | $[64, 512]_{64}$ |
| d_dropout_first | No | TabResNet | $\{0; .5\}_{.05}$ |
| d_dropout_second | No | TabResNet | $\{0; .5\}_{.05}$ |
| batch_norm | No | VGG | $\{\text{True, False}\}$ |
| zero_init_residual | No | ResNet34, ResNet50 | $\{\text{True, False}\}$ |

networks. Following the empirical evaluation in (Franc et al., 2023), we split the training data in half to train SELE and REG networks: on the one half we train the classifier, on the other half, the uncertainty network. We provide the best configurations we employed in the final analysis in Tables A4-A12.

Table A4: Best configurations for DG, divided by architectures.

| Dataset | optim. | l. rate w. | decay t. | decay mom. | SGD nest. | o | n_blocks | d_token | d_att. | drop. | res. | drop. | d_ffn_factor | ffn_drop. | Arch. | |
|---|--------|------------|----------|------------|-----------|-------|----------|---------|--------|-------|------|-------|--------------|-----------|---------------|--|
| adult | SGD | 1e-03 | 1e-04 | False | False | 1.0 | 1 | 320 | .30 | | | | 2.67 | .45 | FTTransformer | |
| compass | Adam | 1e-04 | 1e-04 | False | False | 2.0 | 2 | 64 | .20 | | | | 1.33 | .10 | | |
| covtype | Adam | 1e-04 | 1e-04 | False | False | 4.2 | 1 | 512 | .15 | | | | 1.67 | .50 | | |
| electricity | Adam | 1e-05 | 1e-03 | False | False | 2.0 | 1 | 192 | .05 | | | | 1.00 | .45 | | |
| eye | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| heloc | AdamW | 1e-05 | 1e-05 | True | True | 1.8 | 3 | 384 | .25 | .10 | | | 1.33 | .45 | | |
| higgs | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| house | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| jam1s | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| kddipmsy7 | Adam | 1e-03 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| letter | Adam | 1e-05 | 1e-03 | False | False | 18.4 | 4 | 64 | .05 | | | | 1.00 | .40 | | |
| magic | Adam | 1e-03 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| online | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.67 | .50 | | |
| phoneme | Adam | 1e-03 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| pol | AdamW | 1e-03 | 1e-05 | True | True | 8.4 | 4 | 256 | .05 | .10 | | | 1.67 | .40 | | |
| r1 | Adam | 1e-04 | 1e-03 | True | True | 1.4 | 2 | 320 | .30 | | | | 1.67 | .15 | | |
| upselling | Adam | 1e-05 | 1e-04 | False | False | 2.0 | 1 | 64 | .15 | | | | 1.33 | .45 | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. o n_blocks d_token d_main d_hidden_f d_drop_first d_drop_sec. Arch. | | | | | | | | | | | | | | | | |
| aloi | SGD | 1e-02 | 1e-04 | False | True | 748.4 | 4 | 128 | 192 | 2.00 | | | .20 | .40 | TabResnet | |
| bank | SGD | 1e-01 | 1e-03 | True | False | 1.2 | 1 | 64 | 192 | 1.00 | | | .45 | .30 | | |
| g1veme | AdamW | 1e-05 | 1e-04 | True | True | 2.0 | 1 | 384 | 512 | 4.00 | | | .40 | .30 | | |
| helena | AdamW | 1e-04 | 1e-03 | True | True | 89.0 | 1 | 384 | 448 | 3.00 | | | .40 | .25 | | |
| indian | SGD | 1e-01 | 1e-03 | True | False | 3.2 | 4 | 64 | 192 | 1.00 | | | .40 | .30 | | |
| miniboone | AdamW | 1e-04 | 1e-04 | True | True | 2.0 | 1 | 384 | 448 | 3.00 | | | .40 | .25 | | |
| microcredit | AdamW | 1e-05 | 1e-04 | True | True | 2.0 | 1 | 448 | 448 | 3.00 | | | .30 | .25 | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. o zero_init_resid Arch. | | | | | | | | | | | | | | | | |
| MNIST | | | | | | | | | | | | | | | | |
| bloodm1st | Adam | 1e-04 | 1e-04 | True | True | | | | | 2.8 | | | True | | | |
| breastm1st | Adam | 1e-04 | 1e-04 | True | True | | | | | 4.0 | | | False | | | |
| chestm1st | SGD | 1e-01 | 1e-04 | True | True | | | | | 5.0 | | | False | | | |
| dermam1st | SGD | 1e-01 | 1e-04 | True | True | | | | | 5.0 | | | False | | | |
| food101 | AdamW | 1e-04 | 1e-03 | True | True | | | | | 3.4 | | | True | | | |
| food101 | Adam | 1e-03 | 1e-06 | False | False | | | | | 31.4 | | | True | | | |
| occhn1st | SGD | 1e-01 | 1e-04 | True | True | | | | | 5.0 | | | False | | | |
| organm1st | SGD | 1e-02 | 1e-03 | False | False | | | | | .93 | | | True | | | |
| organm1st | SGD | 1e-02 | 1e-03 | False | False | | | | | 7.2 | | | True | | | |
| organm1st | SGD | 1e-02 | 1e-03 | False | False | | | | | 6.8 | | | True | | | |
| organm1st | SGD | 1e-02 | 1e-03 | False | False | | | | | 5.8 | | | True | | | |
| oxfordp1st | AdamW | 1e-04 | 1e-03 | True | True | | | | | .9 | | | False | 1.2 | | |
| oxfordp1st | AdamW | 1e-04 | 1e-03 | True | True | | | | | 1.2 | | | True | 3.4 | | |
| pathm1st | Adam | 1e-03 | 1e-05 | False | False | | | | | 3.4 | | | True | 5.0 | | |
| pneumoniam1st | SGD | 1e-01 | 1e-04 | True | True | | | | | 5.0 | | | False | 2.8 | | |
| reitm1st | AdamW | 1e-04 | 1e-03 | True | True | | | | | 2.8 | | | True | 8.0 | | |
| stanfordcars | Adam | 1e-04 | 1e-06 | False | False | | | | | 161.0 | | | True | 2.8 | | |
| stanfordcars | Adam | 1e-04 | 1e-06 | False | False | | | | | 8.0 | | | True | 2.8 | | |
| tissuem1st | SGD | 1e-03 | 1e-03 | True | True | | | | | .96 | | | True | 8.0 | | |
| tissuem1st | SGD | 1e-03 | 1e-03 | True | True | | | | | 8.0 | | | True | 8.0 | | |
| waterbirds | AdamW | 1e-04 | 1e-06 | False | False | | | | | 2.0 | | | False | 2.0 | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. o b_norm Arch. | | | | | | | | | | | | | | | | |
| SVHN | | | | | | | | | | | | | | | | |
| catdogs | SGD | 1e-04 | 1e-05 | False | False | | | | | .90 | | | True | 7.8 | True | |
| catdogs | SGD | 1e-02 | 1e-03 | False | False | | | | | .91 | | | True | 1.4 | True | |
| cat10 | SGD | 1e-02 | 1e-03 | False | False | | | | | .91 | | | True | 3.2 | True | |
| Cifar10 | | | | | | | | | | | | | | | | |
| cat10 | SGD | 1e-02 | 1e-03 | False | False | | | | | .91 | | | True | 3.2 | True | |
| Cifar10 | | | | | | | | | | | | | | | | |
| cat10 | SGD | 1e-02 | 1e-03 | False | False | | | | | .91 | | | True | 3.2 | True | |

Table A5: Best configurations for SAT, divided by architectures.

| Dataset | optim. | l. | rate | w. | decay | t. | decay | mom. | SGD | nest. | E_s | γ | n_blocks | d_token | res. | drop. | d_ffn_factor | ffn_drop. | Arch. | | |
|-------------|--------|-------|-------|-------|-------|-----|-------|------|-----|-------|-------|----------|----------|---------|------|-------|--------------|-----------|-------|--|------|
| adult | AdamW | 1e-05 | 1e-05 | True | 45 | .97 | 3 | 128 | .40 | 1.67 | | | | | | | | | | | |
| compass | AdamW | 1e-05 | 1e-05 | False | 45 | .92 | 2 | 128 | .50 | 2.33 | | | | | | | | | | | .45 |
| covtype | Adam | 1e-04 | 1e-06 | True | 15 | .95 | 4 | 320 | .15 | 1.67 | | | | | | | | | | | .405 |
| electricity | Adam | 1e-04 | 1e-06 | False | 45 | .98 | 3 | 128 | .45 | 1.67 | | | | | | | | | | | .10 |
| eye | AdamW | 1e-02 | 1e-06 | True | 45 | .97 | 3 | 256 | .20 | 2.33 | | | | | | | | | | | .45 |
| heloc | AdamW | 1e-05 | 1e-05 | False | 43 | .92 | 2 | 192 | .45 | 2.00 | | | | | | | | | | | |
| higgs | Adam | 1e-04 | 1e-03 | True | 60 | .99 | 2 | 192 | .15 | 1.00 | | | | | | | | | | | .35 |
| house | AdamW | 1e-03 | 1e-06 | False | 19 | .97 | 1 | 384 | .25 | 1.00 | | | | | | | | | | | .20 |
| jannis | Adam | 1e-04 | 1e-06 | True | 19 | .94 | 3 | 128 | .25 | 1.00 | | | | | | | | | | | .30 |
| kddipums97 | Adam | 1e-04 | 1e-06 | True | 19 | .93 | 4 | 128 | .25 | 1.00 | | | | | | | | | | | .30 |
| letter | Adam | 1e-04 | 1e-06 | True | 19 | .92 | 4 | 256 | .25 | 1.00 | | | | | | | | | | | .30 |
| magic | Adam | 1e-03 | 1e-06 | True | 13 | .94 | 4 | 128 | .25 | 1.00 | | | | | | | | | | | .30 |
| online | Adam | 1e-04 | 1e-05 | True | 60 | .99 | 2 | 192 | .15 | 1.00 | | | | | | | | | | | .10 |
| phoneme | Adam | 1e-03 | 1e-06 | True | 19 | .93 | 4 | 64 | .30 | 1.00 | | | | | | | | | | | .30 |
| pol | Adam | 1e-03 | 1e-06 | True | 19 | .93 | 4 | 192 | .15 | 1.00 | | | | | | | | | | | .30 |
| rl | Adam | 1e-04 | 1e-06 | False | 0 | .98 | 3 | 64 | .45 | 1.00 | | | | | | | | | | | .10 |
| upselling | Adam | 1e-04 | 1e-06 | True | 15 | .93 | 4 | 128 | .25 | 1.00 | | | | | | | | | | | .30 |

| Dataset | optim. | l. | rate | w. | decay | t. | decay | mom. | SGD | nest. | E_s | γ | n_blocks | d_token | d_main | d_hidden | f | d_drop | first | d_drop | sec. | Arch. | |
|-----------|--------|-------|-------|-------|-------|-----|-------|------|-----|-------|-------|----------|----------|---------|--------|----------|---|--------|-------|--------|------|-------|--|
| aloi | Adam | 1e-04 | 1e-05 | False | 0 | .93 | 2 | 320 | 384 | 3.00 | .50 | | | | | | | | | | | .50 | |
| bank | Adam | 1e-02 | 1e-04 | True | 15 | .93 | 2 | 320 | 320 | 2.00 | .45 | | | | | | | | | | | .15 | |
| giveme | Adam | 1e-04 | 1e-04 | False | 30 | .91 | 1 | 512 | 192 | 3.00 | .15 | | | | | | | | | | | .25 | |
| helena | Adam | 1e-04 | 1e-04 | False | 30 | .9 | 4 | 192 | 128 | 2.00 | .45 | | | | | | | | | | | .15 | |
| indian | Adam | 1e-03 | 1e-03 | True | 0 | .98 | 4 | 512 | 192 | 1.00 | .10 | | | | | | | | | | | .20 | |
| miniboone | Adam | 1e-03 | 1e-03 | True | 0 | .98 | 4 | 512 | 192 | 1.00 | .10 | | | | | | | | | | | .20 | |
| ucicredit | Adam | 1e-03 | 1e-03 | True | 0 | .98 | 4 | 512 | 192 | 1.00 | .10 | | | | | | | | | | | .20 | |

| Dataset | optim. | l. | rate | w. | decay | t. | decay | mom. | SGD | nest. | E_s | γ | zero_init | resid | Arch. |
|---------------|--------|-------|-------|-------|-------|-------|-------|------|-------|-------|-------|----------|-----------|-------|-------|
| MNIST | SGD | 1e-02 | 1e-05 | False | .92 | False | 0 | .99 | True | True | | | | | |
| bloodmist | Adam | 1e-02 | 1e-06 | True | 60 | .9 | 60 | .9 | True | True | | | | | |
| breastmist | AdamW | 1e-05 | 1e-03 | False | 0 | .91 | 0 | .91 | False | False | | | | | |
| chestmist | Adam | 1e-02 | 1e-06 | True | 60 | .9 | 60 | .9 | True | True | | | | | |
| dermamist | AdamW | 1e-04 | 1e-05 | True | 30 | .94 | 30 | .94 | True | True | | | | | |
| food101 | Adam | 1e-02 | 1e-06 | True | 45 | .9 | 45 | .9 | False | False | | | | | |
| octmnist | AdamW | 1e-03 | 1e-05 | True | 60 | .92 | 60 | .92 | True | True | | | | | |
| organamist | Adam | 1e-03 | 1e-06 | False | 0 | .98 | 0 | .98 | True | True | | | | | |
| organamist | Adam | 1e-03 | 1e-06 | True | 0 | .91 | 0 | .91 | True | True | | | | | |
| organamist | Adam | 1e-03 | 1e-06 | True | 45 | .93 | 45 | .93 | True | True | | | | | |
| oxfordpets | AdamW | 1e-03 | 1e-06 | True | 45 | .91 | 45 | .91 | True | True | | | | | |
| oxfordpets | AdamW | 1e-03 | 1e-06 | False | 15 | .91 | 15 | .91 | False | False | | | | | |
| pathmnist | Adam | 1e-02 | 1e-06 | True | 15 | .91 | 15 | .91 | True | True | | | | | |
| pneumoniamist | SGD | 1e-01 | 1e-04 | True | .9 | False | 60 | .9 | True | True | | | | | |
| retinamist | AdamW | 1e-04 | 1e-06 | True | 15 | .93 | 15 | .93 | True | True | | | | | |
| stanfordcars | Adam | 1e-03 | 1e-05 | False | 15 | .93 | 15 | .93 | True | True | | | | | |
| tissuemist | Adam | 1e-02 | 1e-06 | True | 15 | .91 | 15 | .91 | True | True | | | | | |
| waterbirds | SGD | 1e-02 | 1e-03 | False | .9 | True | 30 | .92 | False | False | | | | | |

| Dataset | optim. | l. | rate | w. | decay | t. | decay | mom. | SGD | nest. | E_s | γ | b_norm | Arch. |
|-----------|--------|-------|-------|------|-------|-----|-------|------|------|-------|-------|----------|--------|-------|
| SVMN | Adam | 1e-02 | 1e-06 | True | 45 | .9 | 45 | .9 | True | True | | | ✓ | |
| catsdogs | Adam | 1e-02 | 1e-06 | True | 45 | .9 | 45 | .9 | True | True | | | ✓ | |
| cifario10 | Adam | 1e-02 | 1e-06 | True | 45 | .93 | 45 | .93 | True | True | | | ✓ | |

Table A6: Best configurations for SAT+EM, divided by architectures.

| Dataset | optim. | l. rate | w. decay | t. decay | mom. SGD | nest. | E_s | γ | β | n_blocks | d_token | atc. | drop. | res. | drop. | d_frn | factor | frn | drop. | Arch. | |
|--|--------|---------|----------|----------|----------|-------|-------|----------|---------|----------|---------|------|-------|------|-------|-------|--------|-----|-------|-------|--|
| adult | SGID | 1e-04 | 1e-06 | False | .95 | True | 45 | .9 | .01 | 3 | 256 | .50 | | | | | | | | .50 | |
| compass | Adam | 1e-03 | 1e-06 | False | | | 60 | .96 | .0001 | 3 | 512 | .45 | | | | | | | | 2.67 | |
| covertype | AdamW | 1e-05 | 1e-06 | False | | | 30 | .98 | .001 | 3 | 320 | .10 | | | | | | | | 1.33 | |
| electricity | AdamW | 1e-05 | 1e-05 | False | | | 30 | .96 | .001 | 3 | 320 | .05 | | | | | | | | 2.33 | |
| eye | AdamW | 1e-05 | 1e-04 | True | | | 30 | .91 | .01 | 2 | 448 | .15 | | | | | | | | 1.33 | |
| heloc | AdamW | 1e-03 | 1e-05 | True | | | 30 | .92 | .0001 | 2 | 384 | .10 | | | | | | | | 2.33 | |
| higgs | AdamW | 1e-03 | 1e-05 | True | | | 30 | .98 | .0001 | 2 | 448 | .10 | | | | | | | | 1.67 | |
| house | Adam | 1e-03 | 1e-06 | False | | | 60 | .96 | .0001 | 3 | 512 | .45 | | | | | | | | 2.00 | |
| jannis | AdamW | 1e-03 | 1e-06 | True | | | 30 | .97 | .0001 | 3 | 448 | .10 | | | | | | | | 2.00 | |
| Kddipums97 | Adam | 1e-03 | 1e-06 | False | | | 60 | .95 | .0001 | 3 | 512 | .45 | | | | | | | | 2.67 | |
| letter | AdamW | 1e-05 | 1e-05 | False | | | 30 | .96 | .001 | 3 | 384 | .10 | | | | | | | | 1.33 | |
| magic | Adam | 1e-03 | 1e-06 | False | | | 60 | .94 | .0001 | 1 | 512 | .45 | | | | | | | | 1.33 | |
| online | Adam | 1e-04 | 1e-06 | False | | | 15 | .99 | .01 | 3 | 320 | .20 | | | | | | | | 2.00 | |
| phoneme | AdamW | 1e-04 | 1e-06 | True | | | 30 | .98 | .0001 | 3 | 448 | .15 | | | | | | | | 2.00 | |
| pol | AdamW | 1e-04 | 1e-06 | True | | | 60 | .97 | .0001 | 1 | 512 | .50 | | | | | | | | 1.33 | |
| rl | Adam | 1e-04 | 1e-06 | False | | | 30 | .99 | .0001 | 3 | 448 | .10 | | | | | | | | 2.00 | |
| upselling | AdamW | 1e-05 | 1e-06 | True | | | 30 | .99 | .0001 | 1 | 512 | .10 | | | | | | | | 2.00 | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. E_s γ β n_blocks d_token d_main d_hidden_f d_drop_first d_drop_sec. Arch. | | | | | | | | | | | | | | | | | | | | | |
| aloi | SGID | 1e-01 | 1e-04 | False | .93 | False | 15 | .91 | .001 | 3 | 192 | 512 | 4.00 | | | | | | | .10 | |
| bank | Adam | 1e-02 | 1e-06 | True | | | 15 | .94 | .001 | 1 | 384 | 256 | 1.00 | | | | | | | .30 | |
| giveme | Adam | 1e-04 | 1e-06 | True | | | 0 | .93 | .0001 | 2 | 384 | 256 | 2.00 | | | | | | | .25 | |
| helena | Adam | 1e-04 | 1e-06 | True | | | 15 | .95 | .0001 | 2 | 384 | 256 | 1.00 | | | | | | | .05 | |
| indian | Adam | 1e-03 | 1e-06 | True | | | 0 | .94 | .001 | 2 | 384 | 256 | 1.00 | | | | | | | .30 | |
| miniboone | Adam | 1e-02 | 1e-06 | True | | | 15 | .94 | .001 | 1 | 384 | 256 | 1.00 | | | | | | | .30 | |
| notcredit | AdamW | 1e-05 | 1e-06 | False | | | 60 | .92 | .001 | 2 | 512 | 128 | 3.00 | | | | | | | .40 | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. E_s γ β zero_init_resid Arch. | | | | | | | | | | | | | | | | | | | | | |
| MNIST | SGID | 1e-01 | 1e-05 | False | .96 | True | 60 | .98 | .0001 | True | | | True | | | | | | | | |
| bloodmst | Adam | 1e-03 | 1e-06 | False | | | 60 | .96 | .0001 | False | | | False | | | | | | | | |
| breastmst | SGID | 1e-02 | 1e-04 | True | .9 | False | 45 | .9 | .01 | False | | | False | | | | | | | | |
| chestmst | Adam | 1e-03 | 1e-06 | True | | | 0 | .93 | .0001 | False | | | False | | | | | | | | |
| dermst | SGID | 1e-04 | 1e-05 | False | .99 | True | 60 | .98 | .0001 | True | | | True | | | | | | | | |
| food101 | Adam | 1e-02 | 1e-06 | True | | | 0 | .9 | .01 | True | | | True | | | | | | | | |
| octmst | Adam | 1e-03 | 1e-06 | True | | | 0 | .93 | .0001 | False | | | False | | | | | | | | |
| organmst | Adam | 1e-04 | 1e-06 | False | | | 0 | .96 | .0001 | False | | | False | | | | | | | | |
| organmst | Adam | 1e-03 | 1e-06 | False | | | 60 | .98 | .0001 | False | | | False | | | | | | | | |
| organmst | SGID | 1e-02 | 1e-05 | False | .95 | True | 60 | .98 | .0001 | True | | | True | | | | | | | | |
| oxfordpets | SGID | 1e-02 | 1e-04 | True | .93 | True | 0 | .98 | .0001 | False | | | False | | | | | | | | |
| pathmst | Adam | 1e-03 | 1e-06 | False | | | 0 | .95 | .0001 | False | | | False | | | | | | | | |
| pneumoniamst | Adam | 1e-05 | 1e-06 | False | | | 0 | .94 | .0001 | False | | | False | | | | | | | | |
| retiamst | Adam | 1e-04 | 1e-06 | True | | | 0 | .93 | .0001 | False | | | False | | | | | | | | |
| stanfordcars | Adam | 1e-03 | 1e-04 | False | | | 15 | .98 | .001 | True | | | True | | | | | | | | |
| tissuemst | Adam | 1e-02 | 1e-06 | True | | | 0 | .93 | .0001 | False | | | False | | | | | | | | |
| waterbirds | Adam | 1e-03 | 1e-06 | False | | | 0 | .96 | .0001 | False | | | False | | | | | | | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. E_s γ β b_norm Arch. | | | | | | | | | | | | | | | | | | | | | |
| SVHN | Adam | 1e-04 | 1e-06 | True | .96 | True | 0 | .93 | .0001 | False | | | False | | | | | | | | |
| catsdogs | SGID | 1e-01 | 1e-05 | False | .93 | True | 60 | .98 | .0001 | True | | | True | | | | | | | | |
| char10 | SGID | 1e-01 | 1e-04 | False | | | 60 | .98 | .0001 | True | | | True | | | | | | | | |
| Resnet18-50 | | | | | | | | | | | | | | | | | | | | | |
| FTTtransformer | | | | | | | | | | | | | | | | | | | | | |
| TabResnet | | | | | | | | | | | | | | | | | | | | | |
| VG | | | | | | | | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | | | | | | | | |

Table A7: Best configurations for SELNET, divided by architectures.

| Dataset | optim. | l. | rate | w. | decay | t. | decay | mom. | SGD | nest. | λ | α | n_blocks | d_token | att. | d_drop. | res. | d_ffn_factor | ffn_drop. | Arch. | |
|--|--------|-------|-------|-------|-------|----|-------|------|-----|-------|-----------|----------|----------|---------|------|---------|------|--------------|-----------|-------|--|
| adult | AdamW | 1e-05 | 1e-06 | True | | | | | | | 32.65 | | 3 | 192 | .40 | | | | 2.00 | | |
| compass | Adam | 1e-04 | 1e-05 | True | | | | | | | 64.75 | | 2 | 192 | .15 | | .10 | | 1.00 | | |
| covtype | Adam | 1e-04 | 1e-06 | True | | | | | | | 10.25 | | 3 | 384 | .10 | | | | .67 | | |
| electricity | Adam | 1e-04 | 1e-04 | True | | | | | | | 64.75 | | 2 | 320 | .35 | | .05 | | 1.00 | | |
| eye | AdamW | 1e-05 | 1e-06 | True | | | | | | | 32.70 | | 3 | 320 | .35 | | | | 2.67 | .05 | |
| heLoc | Adam | 1e-05 | 1e-06 | False | | | | | | | 8.70 | | 3 | 64 | .50 | | | | 1.33 | .10 | |
| higgs | AdamW | 1e-04 | 1e-06 | True | | | | | | | 64.65 | | 3 | 192 | .25 | | | | 2.67 | | |
| house | Adam | 1e-04 | 1e-05 | True | | | | | | | 64.75 | | 2 | 192 | .15 | | .10 | | 1.00 | .50 | |
| jaannis | AdamW | 1e-04 | 1e-05 | False | | | | | | | 64.35 | | 2 | 256 | .50 | | | | 2.33 | | |
| kddipums97 | Adam | 1e-04 | 1e-06 | False | | | | | | | 8.70 | | 3 | 64 | .50 | | | | 1.33 | .10 | |
| Letter | Adam | 1e-04 | 1e-06 | True | | | | | | | 16.40 | | 4 | 320 | .10 | | | | .67 | .20 | |
| magic | Adam | 1e-04 | 1e-06 | True | | | | | | | 32.55 | | 3 | 256 | .15 | | .20 | | .67 | .10 | |
| online | Adam | 1e-03 | 1e-05 | True | | | | | | | 64.70 | | 2 | 192 | .15 | | | | 1.00 | .30 | |
| phoneme | Adam | 1e-03 | 1e-06 | True | | | | | | | 8.35 | | 3 | 320 | .10 | | | | 1.33 | | |
| pol | Adam | 1e-04 | 1e-05 | True | | | | | | | 64.75 | | 2 | 192 | .15 | | .10 | | 1.00 | | |
| rl | Adam | 1e-03 | 1e-06 | False | | | | | | | 8.70 | | 3 | 64 | .50 | | | | 1.00 | .10 | |
| upselling | Adam | 1e-04 | 1e-06 | True | | | | | | | 16.45 | | 4 | 128 | .25 | | | | 1.00 | .30 | |
| FTTransformer | | | | | | | | | | | | | | | | | | | | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. λ α n_blocks d_token d_main d_hidden_f d_drop_first d_drop_sec. Arch. | | | | | | | | | | | | | | | | | | | | | |
| aloi | Adam | 1e-04 | 1e-05 | False | | | | | | | 32.25 | | 4 | 128 | 128 | 2.00 | | | .50 | .05 | |
| bank | Adam | 1e-03 | 1e-03 | True | | | | | | | 8.70 | | 4 | 512 | 192 | 1.00 | | | .10 | .20 | |
| giveme | Adam | 1e-03 | 1e-03 | True | | | | | | | 8.70 | | 4 | 512 | 192 | 1.00 | | | .10 | .20 | |
| helena | Adam | 1e-04 | 1e-04 | False | | | | | | | 32.25 | | 4 | 128 | 128 | 2.00 | | | .50 | .10 | |
| indian | AdamW | 1e-02 | 1e-04 | True | | | | | | | 32.65 | | 2 | 256 | 320 | 2.00 | | | .40 | .15 | |
| miniboone | Adam | 1e-03 | 1e-03 | True | | | | | | | 8.70 | | 4 | 512 | 192 | 1.00 | | | .10 | .20 | |
| ucicredit | Adam | 1e-04 | 1e-04 | False | | | | | | | 32.25 | | 4 | 192 | 64 | 2.00 | | | .45 | .10 | |
| TabResnet | | | | | | | | | | | | | | | | | | | | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. λ α zero_init_resid Arch. | | | | | | | | | | | | | | | | | | | | | |
| MNIST | SGD | 1e-01 | 1e-04 | True | | | | | | | .94 | | True | 32 | .55 | | | | False | | |
| bloodmst | SGD | 1e-01 | 1e-04 | True | | | | | | | .99 | | True | 32 | .45 | | | | False | | |
| breastmst | AdamW | 1e-03 | 1e-05 | True | | | | | | | 32.45 | | True | 32 | .45 | | | | True | | |
| chestmst | Adam | 1e-02 | 1e-06 | True | | | | | | | 64.25 | | True | 64 | .25 | | | | True | | |
| dermamst | AdamW | 1e-05 | 1e-03 | False | | | | | | | 16.30 | | False | 16 | .30 | | | | False | | |
| food101 | SGD | 1e-01 | 1e-04 | True | | | | | | | .9 | | False | 32 | .50 | | | | True | | |
| octmst | Adam | 1e-02 | 1e-06 | True | | | | | | | .94 | | True | 16 | .30 | | | | True | | |
| organamnist | SGD | 1e-01 | 1e-04 | True | | | | | | | .94 | | True | 32 | .50 | | | | True | | |
| organamnist | Adam | 1e-03 | 1e-06 | True | | | | | | | 16.25 | | True | 16 | .25 | | | | True | | |
| organamnist | Adam | 1e-03 | 1e-06 | True | | | | | | | 16.25 | | True | 16 | .25 | | | | True | | |
| oxfordpets | AdamW | 1e-03 | 1e-06 | False | | | | | | | .93 | | True | 64 | .60 | | | | False | | |
| pathmst | SGD | 1e-01 | 1e-04 | True | | | | | | | 16.30 | | True | 16 | .30 | | | | True | | |
| pneumiamnist | Adam | 1e-02 | 1e-06 | True | | | | | | | 64.60 | | True | 64 | .60 | | | | True | | |
| retinamnist | Adam | 1e-02 | 1e-06 | True | | | | | | | 16.30 | | True | 16 | .30 | | | | True | | |
| stanfordcars | Adam | 1e-03 | 1e-06 | True | | | | | | | 64.25 | | True | 64 | .25 | | | | True | | |
| tissuemnist | Adam | 1e-02 | 1e-06 | True | | | | | | | 16.30 | | True | 16 | .30 | | | | True | | |
| waterbirds | Adam | 1e-03 | 1e-06 | False | | | | | | | 8.45 | | True | 8 | .45 | | | | True | | |
| Resnet18-50 | | | | | | | | | | | | | | | | | | | | | |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. λ α b_norm Arch. | | | | | | | | | | | | | | | | | | | | | |
| SVMN | Adam | 1e-02 | 1e-06 | True | | | | | | | 64.25 | | True | 64 | .25 | | | | True | V | |
| catsdogs | Adam | 1e-02 | 1e-06 | True | | | | | | | 32.25 | | True | 32 | .25 | | | | True | C | |
| cifar10 | Adam | 1e-03 | 1e-06 | False | | | | | | | 8.70 | | True | 8 | .70 | | | | True | C | |

Table A8: Best configurations for SELNET+EM, divided by architectures.

| Dataset | optim. | l. rate | w. decay | t. decay | mon. SGD | nest. | λ | α | β | n_blocks | d_token | attn. | drop. | res. | drop. | d_ffn | factor | ffn | drop. | Arch. |
|---------------|--------|---------|----------|----------|----------|-------|-----------|----------|---------|----------|---------|-------|----------|------|-------|-------|--------|------|-------|-------|
| adult | AdamW | 1e-02 | 1e-05 | True | False | 16 | .25 | .0001 | 2 | 384 | .15 | 2.33 | .50 | | | | | | | |
| compass | AdamW | 1e-05 | 1e-06 | True | True | 16 | .75 | .0001 | 3 | 448 | .15 | 2.00 | .30 | | | | | | | |
| covertype | AdamW | 1e-05 | 1e-05 | False | False | 16 | .60 | .001 | 3 | 320 | .05 | 2.33 | .10 | | | | | | | |
| electricity | AdamW | 1e-05 | 1e-06 | False | False | 16 | .45 | .0001 | 4 | 448 | .10 | 1.00 | .35 | | | | | | | |
| eye | AdamW | 1e-04 | 1e-04 | True | True | 64 | .25 | .01 | 1 | 448 | .25 | 1.67 | .20 | | | | | | | |
| heloc | AdamW | 1e-03 | 1e-04 | True | True | 64 | .25 | .01 | 1 | 448 | .35 | 1.33 | .15 | | | | | | | |
| higgs | AdamW | 1e-03 | 1e-04 | True | True | 64 | .35 | .01 | 1 | 384 | .20 | 1.33 | .10 | | | | | | | |
| house | AdamW | 1e-05 | 1e-05 | True | True | 8 | .80 | .0001 | 3 | 384 | .50 | 1.33 | .10 | | | | | | | |
| jamnis | Adam | 1e-04 | 1e-06 | False | False | 64 | .60 | .0001 | 1 | 512 | .20 | 1.67 | .15 | | | | | | | |
| kddipwm57 | AdamW | 1e-03 | 1e-06 | True | True | 16 | .70 | .0001 | 3 | 448 | .10 | 2.00 | .30 | | | | | | | |
| letter | AdamW | 1e-04 | 1e-06 | True | True | 16 | .70 | .0001 | 3 | 448 | .10 | 2.00 | .30 | | | | | | | |
| magic | AdamW | 1e-03 | 1e-06 | True | True | 16 | .75 | .0001 | 3 | 448 | .15 | 2.00 | .25 | | | | | | | |
| online | AdamW | 1e-03 | 1e-06 | True | True | 16 | .75 | .0001 | 3 | 448 | .15 | 2.00 | .25 | | | | | | | |
| phoneme | AdamW | 1e-05 | 1e-06 | False | False | 16 | .30 | .0001 | 4 | 192 | .15 | 1.00 | .30 | | | | | | | |
| pol | AdamW | 1e-05 | 1e-05 | False | False | 16 | .65 | .001 | 3 | 320 | .20 | 1.67 | .35 | | | | | | | |
| ri | Adam | 1e-04 | 1e-06 | False | False | 16 | .55 | .0001 | 1 | 512 | .45 | 2.00 | .10 | | | | | | | |
| upspelling | AdamW | 1e-05 | 1e-06 | True | True | 16 | .70 | .0001 | 3 | 448 | .10 | 2.00 | .30 | | | | | | | |
| Dataset | optim. | l. rate | w. decay | t. decay | mon. SGD | nest. | λ | α | β | n_blocks | d_token | main | d_hidden | f_d | drop | first | d_drop | sec. | Arch. | |
| aloi | SGD | 1e-01 | 1e-05 | False | False | 64 | .65 | .0001 | 4 | 320 | 384 | 3.00 | .15 | | | | | | | |
| bank | Adam | 1e-03 | 1e-06 | True | True | 16 | .55 | .0001 | 1 | 320 | 320 | 1.00 | .30 | | | | | | | |
| giveme | Adam | 1e-03 | 1e-06 | True | True | 16 | .45 | .0001 | 1 | 384 | 256 | 1.00 | .25 | | | | | | | |
| helena | Adam | 1e-03 | 1e-06 | True | True | 8 | .50 | .0001 | 2 | 448 | 192 | 1.00 | .40 | | | | | | | |
| indian | AdamW | 1e-05 | 1e-06 | False | False | 64 | .35 | .001 | 2 | 512 | 128 | 3.00 | .40 | | | | | | | |
| miniboone | Adam | 1e-04 | 1e-06 | True | True | 8 | .45 | .0001 | 2 | 384 | 256 | 2.00 | .25 | | | | | | | |
| netcredit | Adam | 1e-03 | 1e-06 | True | True | 16 | .45 | .001 | 1 | 384 | 256 | 1.00 | .30 | | | | | | | |
| Dataset | optim. | l. rate | w. decay | t. decay | mon. SGD | nest. | λ | α | β | zero | init | resid | Arch. | | | | | | | |
| MNIST | SGD | 1e-03 | 1e-05 | False | False | .93 | False | 64 | .35 | .0001 | False | False | Resnet18 | 50 | | | | | | |
| bloodmst | Adam | 1e-03 | 1e-06 | False | False | 64 | .55 | .0001 | False | False | False | False | Arch. | | | | | | | |
| breastmst | Adam | 1e-04 | 1e-06 | True | True | 8 | .40 | .0001 | False | False | False | False | Arch. | | | | | | | |
| chestmst | SGD | 1e-02 | 1e-05 | False | False | .93 | True | 64 | .30 | .0001 | False | False | Arch. | | | | | | | |
| dermmst | Adam | 1e-05 | 1e-06 | False | False | 16 | .60 | .0001 | False | False | False | False | Arch. | | | | | | | |
| food101 | Adam | 1e-04 | 1e-06 | True | True | 8 | .40 | .0001 | False | False | False | False | Arch. | | | | | | | |
| ocnumst | Adam | 1e-04 | 1e-06 | True | True | 16 | .60 | .0001 | False | False | False | False | Arch. | | | | | | | |
| organmst | Adam | 1e-03 | 1e-06 | True | True | 8 | .55 | .0001 | False | False | False | False | Arch. | | | | | | | |
| organmst | Adam | 1e-04 | 1e-06 | True | True | 64 | .55 | .0001 | False | False | False | False | Arch. | | | | | | | |
| organmst | Adam | 1e-04 | 1e-06 | True | True | 64 | .60 | .0001 | False | False | False | False | Arch. | | | | | | | |
| oxfordpets | SGD | 1e-02 | 1e-04 | True | True | .95 | False | 32 | .50 | .0001 | False | False | Resnet18 | 50 | | | | | | |
| pathmst | Adam | 1e-03 | 1e-06 | False | False | 64 | .55 | .0001 | False | False | False | False | Arch. | | | | | | | |
| pneumoniamst | Adam | 1e-03 | 1e-06 | False | False | 64 | .55 | .0001 | False | False | False | False | Arch. | | | | | | | |
| retinamst | Adam | 1e-02 | 1e-06 | False | False | 8 | .60 | .0001 | False | False | False | False | Arch. | | | | | | | |
| standfordcars | Adam | 1e-04 | 1e-06 | True | True | 8 | .40 | .0001 | False | False | False | False | Arch. | | | | | | | |
| tissamst | Adam | 1e-04 | 1e-06 | True | True | 8 | .40 | .0001 | False | False | False | False | Arch. | | | | | | | |
| waterbirds | Adam | 1e-04 | 1e-06 | False | False | 16 | .60 | .0001 | False | False | False | False | Arch. | | | | | | | |
| Dataset | optim. | l. rate | w. decay | t. decay | mon. SGD | nest. | λ | α | β | b_norm | Arch. | | | | | | | | | |
| SVHN | Adam | 1e-04 | 1e-06 | True | True | 8 | .40 | .0001 | False | True | Arch. | | | | | | | | | |
| catsdogs | Adam | 1e-03 | 1e-06 | False | False | 16 | .65 | .0001 | False | True | Arch. | | | | | | | | | |
| cifar10 | SGD | 1e-02 | 1e-03 | True | True | .99 | True | 32 | .50 | .0001 | True | True | Arch. | | | | | | | |

Table A10: Best configurations for REG, divided by architectures.

| Dataset | optim. | l. rate v. | decay t. | decay mom. | SGD nest. | optim. | unc. l. | rate unc. | mom. | SGD unc. | nest. | unc. n. | blocks | d.token | att. | drop. | res. | drop. | d.fth. | factor | fth | drop. | Arch. | |
|---|---|------------|----------|------------|-----------|--------|---------|-----------|-------|----------|-------|---------|--------|---------|------|-------|------|-------|--------|--------|-----|-------|-------|-----|
| FTTtransformer | adult | AdamW | 1e-05 | 1e-04 | False | SGD | | 1e-05 | | | | | 2 | 128 | .50 | | | | | | | 2.33 | .45 | |
| | compass | SGD | 1e-01 | 1e-06 | True | AdamW | | 1e-07 | | | | | 4 | 256 | 1.10 | | | | | | | | 1.67 | .35 |
| | covertype | Adam | 1e-03 | 1e-05 | True | AdamW | | 1e-06 | | | | | 2 | 320 | 1.40 | | | | | | | | .67 | .35 |
| | electricity | SGD | 1e-04 | 1e-04 | True | AdamW | | 1e-04 | | | | | 2 | 128 | .35 | | | | | | | | 1.33 | .15 |
| | eye | SGD | 1e-04 | 1e-03 | False | AdamW | | 1e-05 | | | | | 3 | 128 | .35 | | | | | | | | .67 | .25 |
| | heloc | Adam | 1e-03 | 1e-05 | True | SGD | | 1e-04 | | | | | 1 | 64 | 1.15 | | | | | | | | 2.00 | .50 |
| | higgs | AdamW | 1e-04 | 1e-04 | False | SGD | | 1e-04 | | .9 | | | 2 | 64 | .50 | | | | | | | | 2.33 | .25 |
| | house | SGD | 1e-03 | 1e-04 | False | AdamW | | 1e-06 | | | | | 3 | 256 | .35 | | | | | | | | .67 | .15 |
| | jam1s | Adam | 1e-04 | 1e-05 | True | AdamW | | 1e-06 | | | | | 4 | 64 | .25 | | | | | | | | 2.67 | .05 |
| | kddipuns97 | AdamW | 1e-03 | 1e-03 | False | Adam | | 1e-07 | | .9 | | | 1 | 512 | .40 | | | | | | | | 2.67 | .30 |
| | letter | SGD | 1e-02 | 1e-04 | False | AdamW | | 1e-06 | | | | | 3 | 320 | .20 | | | | | | | | 1.67 | .25 |
| | magic | Adam | 1e-04 | 1e-05 | True | AdamW | | 1e-04 | | | | | 1 | 128 | .05 | | | | | | | | 2.67 | .25 |
| | online | Adam | 1e-04 | 1e-06 | False | Adam | | 1e-04 | | | | | 1 | 384 | .30 | | | | | | | | 2.67 | .20 |
| | phone | AdamW | 1e-03 | 1e-03 | False | Adam | | 1e-05 | | | | | 1 | 512 | .40 | | | | | | | | 2.67 | .30 |
| | pol | Adam | 1e-04 | 1e-03 | True | AdamW | | 1e-06 | | | | | 4 | 320 | .25 | | | | | | | | 2.00 | .05 |
| | rs | AdamW | 1e-03 | 1e-03 | False | AdamW | | 1e-06 | | | | | 2 | 128 | .10 | | | | | | | | .67 | .30 |
| | upselling | SGD | 1e-04 | 1e-03 | False | AdamW | | 1e-05 | | .96 | | | 3 | 128 | .35 | | | | | | | | .67 | .15 |
| | Dataset optim. l. rate v. decay t. decay mom. SGD nest. optim. unc. l. rate unc. mom. SGD unc. nest. unc. n.blocks d.token d.main d.hidden_f d.drop_first d.drop_sec. Arch. | | | | | | | | | | | | | | | | | | | | | | | |
| | TabResNet | aloi | Adam | 1e-03 | 1e-06 | True | AdamW | | 1e-06 | | | | 3 | 128 | 512 | 4.00 | | | | | | | .50 | .45 |
| bank | | SGD | 1e-01 | 1e-06 | True | AdamW | | 1e-07 | | | | 4 | 256 | 64 | 3.00 | | | | | | | | .30 | .30 |
| giveme | | Adam | 1e-04 | 1e-05 | False | AdamW | | 1e-04 | | | | 4 | 512 | 192 | 2.00 | | | | | | | | .15 | .15 |
| hahana | | Adam | 1e-04 | 1e-03 | True | SGD | | 1e-06 | | .99 | | | 3 | 192 | 3.00 | | | | | | | | .25 | .05 |
| indian | | AdamW | 1e-04 | 1e-04 | True | Adam | | 1e-03 | | | | | 3 | 256 | 512 | 1.00 | | | | | | | .10 | .05 |
| miniboone | | Adam | 1e-04 | 1e-03 | SGD | Adam | | 1e-05 | | .97 | | | 3 | 256 | 384 | 3.00 | | | | | | | .25 | .15 |
| netcredit | | AdamW | 1e-05 | 1e-04 | True | Adam | | 1e-07 | | | | | 3 | 192 | 512 | 1.00 | | | | | | | .15 | |
| Dataset optim. l. rate v. decay t. decay mom. SGD nest. optim. unc. l. rate unc. mom. SGD unc. nest. unc. zero_init_resid Arch. | | | | | | | | | | | | | | | | | | | | | | | | |
| ResNet18-50 | | mnist | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-07 | | | | | | | | | | | | | | | |
| | | bloodmst | Adam | 1e-02 | 1e-06 | True | SGD | | 1e-03 | | | | | | | | | | | | | | | |
| | | breastmst | Adam | 1e-04 | 1e-04 | True | AdamW | | 1e-06 | | | | | | | | | | | | | | | |
| | | chestrnmt | AdamW | 1e-01 | 1e-04 | True | Adam | | 1e-07 | | | | | | | | | | | | | | | |
| | | dermnmt | SGD | 1e-01 | 1e-04 | True | Adam | | 1e-06 | | .99 | | | | | | | | | | | | | |
| | | food01 | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-06 | | | | | | | | | | | | | | | |
| | | food01t | SGD | 1e-02 | 1e-05 | True | Adam | | 1e-05 | | .99 | | | | | | | | | | | | | |
| | | food01s | SGD | 1e-03 | 1e-06 | False | SGD | | 1e-05 | | | | | | | | | | | | | | | |
| | | organmst | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-07 | | | | | | | | | | | | | | | |
| | | organmst | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-04 | | | | | | | | | | | | | | | |
| | | organmst | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-07 | | | | | | | | | | | | | | | |
| | | organmst | Adam | 1e-04 | 1e-06 | True | Adam | | 1e-04 | | .97 | | | | | | | | | | | | | |
| | organmst | SGD | 1e-01 | 1e-05 | True | Adam | | 1e-07 | | | | | | | | | | | | | | | | |
| | organmst | Adam | 1e-04 | 1e-06 | True | Adam | | 1e-04 | | .97 | | | | | | | | | | | | | | |
| | organmst | SGD | 1e-01 | 1e-06 | True | Adam | | 1e-07 | | | | | | | | | | | | | | | | |
| | organmst | Adam | 1e-03 | 1e-04 | False | Adam | | 1e-03 | | .99 | | | | | | | | | | | | | | |
| | organmst | SGD | 1e-01 | 1e-04 | True | Adam | | 1e-03 | | .99 | | | | | | | | | | | | | | |
| | organmst | Adam | 1e-03 | 1e-04 | False | Adam | | 1e-03 | | .97 | | | | | | | | | | | | | | |
| | organmst | SGD | 1e-01 | 1e-04 | True | Adam | | 1e-06 | | | | | | | | | | | | | | | | |
| | organmst | Adam | 1e-03 | 1e-06 | False | SGD | | 1e-03 | | .97 | | | | | | | | | | | | | | |
| | Dataset optim. l. rate v. decay t. decay mom. SGD nest. optim. unc. l. rate unc. mom. SGD unc. nest. unc. b_norm Arch. | | | | | | | | | | | | | | | | | | | | | | | |
| VGG | svhn | Adam | 1e-04 | 1e-03 | False | AdamW | | 1e-06 | | | | | | | | | | | | | | | | |
| | cars101 | Adam | 1e-02 | 1e-05 | False | AdamW | | 1e-04 | | | | | | | | | | | | | | | | |
| | cars101 | AdamW | 1e-03 | 1e-05 | False | SGD | | 1e-03 | | .99 | | | | | | | | | | | | | | |

Table A11: Best configurations for SELE, divided by architectures.

| Dataset | optim. | 1. rate | v. decay | t. decay | mon. SGD | nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. zero_init_resid | Arch. | |
|---------------|-------------|------------|----------|----------|----------|--------------|--------------|---------------|-------------------|-------------------|---------------|-------------------|----------------------|----------------------|-------|
| FTTransformer | adult | SGD | 1e-03 | 1e-06 | True | AdamW | 1e-07 | .93 | True | AdamW | 1e-07 | .94 | .10 | 1.67 | |
| | compass | SGD | 1e-04 | 1e-03 | False | AdamW | 1e-05 | .96 | True | AdamW | 1e-05 | .94 | .05 | 1.67 | |
| | covtype | SGD | 1e-01 | 1e-06 | True | AdamW | 1e-07 | .9 | False | AdamW | 1e-07 | .94 | .35 | 1.67 | |
| | electricity | Adam | 1e-03 | 1e-06 | False | Adam | 1e-06 | .96 | True | AdamW | 1e-05 | .94 | .15 | 1.67 | |
| | eye | SGD | 1e-03 | 1e-03 | False | AdamW | 1e-03 | .96 | True | AdamW | 1e-03 | .94 | .35 | 1.67 | |
| | heloc | Adam | 1e-03 | 1e-04 | True | SGD | 1e-04 | .96 | True | AdamW | 1e-04 | .94 | .35 | 1.67 | |
| | higgs | AdamW | 1e-04 | 1e-05 | True | AdamW | 1e-04 | .96 | True | AdamW | 1e-04 | .94 | .15 | 2.00 | |
| | house | SGD | 1e-04 | 1e-03 | False | AdamW | 1e-05 | .96 | True | AdamW | 1e-05 | .94 | .15 | .67 | |
| | jannis | SGD | 1e-03 | 1e-03 | False | AdamW | 1e-07 | .94 | True | AdamW | 1e-07 | .94 | .35 | .67 | |
| | house | SGD | 1e-04 | 1e-03 | False | AdamW | 1e-05 | .96 | True | AdamW | 1e-05 | .94 | .35 | .67 | |
| | kdipms97 | SGD | 1e-02 | 1e-03 | False | AdamW | 1e-04 | .95 | True | AdamW | 1e-04 | .99 | .20 | 2.33 | |
| | letter | SGD | 1e-02 | 1e-05 | False | AdamW | 1e-04 | .95 | True | AdamW | 1e-04 | .99 | .20 | 2.33 | |
| | magic | Adam | 1e-03 | 1e-05 | True | Adam | 1e-04 | .95 | True | Adam | 1e-04 | .93 | .15 | .67 | |
| | online | AdamW | 1e-03 | 1e-04 | False | Adam | 1e-04 | .95 | True | Adam | 1e-04 | .93 | .15 | .67 | |
| | phoneme | Adam | 1e-05 | 1e-06 | False | Adam | 1e-04 | .95 | True | Adam | 1e-04 | .93 | .15 | .67 | |
| | pol | Adam | 1e-03 | 1e-05 | True | SGD | 1e-04 | .95 | True | Adam | 1e-04 | .93 | .15 | .67 | |
| r1 | Adam | 1e-05 | 1e-06 | False | Adam | 1e-04 | .95 | True | Adam | 1e-04 | .93 | .15 | .67 | | |
| upselling | Adam | 1e-04 | 1e-05 | True | AdamW | 1e-04 | .95 | True | AdamW | 1e-04 | .93 | .15 | .67 | | |
| Dataset | optim. | 1. rate | v. decay | t. decay | mon. SGD | nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. zero_init_resid | Arch. | |
| Resnet18-50 | aloi | AdamW | 1e-04 | 1e-05 | False | SGD | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | bank | Adam | 1e-02 | 1e-03 | True | Adam | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | givene | AdamW | 1e-04 | 1e-04 | True | Adam | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | helena | Adam | 1e-04 | 1e-06 | True | AdamW | 1e-05 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | indian | Adam | 1e-02 | 1e-03 | True | Adam | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | miniboone | AdamW | 1e-04 | 1e-04 | True | Adam | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | ucicredit | Adam | 1e-02 | 1e-03 | True | Adam | 1e-03 | .92 | True | Adam | 1e-06 | .94 | 3.00 | .40 | |
| | Dataset | optim. | 1. rate | v. decay | t. decay | mon. SGD | nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. zero_init_resid | Arch. |
| | TabResnet | MNIST | SGD | 1e-01 | 1e-03 | True | SGD | 1e-05 | .95 | False | Adam | 1e-06 | .94 | True | True |
| | | bloodmst | Adam | 1e-02 | 1e-06 | False | Adam | 1e-03 | .95 | False | SGD | 1e-04 | .94 | True | True |
| | | breastmst | Adam | 1e-04 | 1e-06 | True | Adam | 1e-03 | .95 | False | AdamW | 1e-04 | .94 | True | True |
| | | chestmst | SGD | 1e-01 | 1e-05 | True | Adam | 1e-03 | .95 | False | Adam | 1e-03 | .94 | True | True |
| | | dermamst | SGD | 1e-01 | 1e-04 | True | Adam | 1e-04 | .93 | False | Adam | 1e-04 | .94 | True | True |
| | | food101 | Adam | 1e-03 | 1e-04 | False | AdamW | 1e-04 | .96 | True | AdamW | 1e-04 | .94 | True | False |
| | | octamst | SGD | 1e-02 | 1e-04 | True | Adam | 1e-04 | .94 | True | Adam | 1e-07 | .94 | True | True |
| | | organeamst | SGD | 1e-01 | 1e-04 | True | Adam | 1e-03 | .94 | True | AdamW | 1e-07 | .94 | True | True |
| organcmst | | Adam | 1e-04 | 1e-03 | False | Adam | 1e-03 | .94 | True | AdamW | 1e-07 | .94 | True | True | |
| organsmst | | Adam | 1e-04 | 1e-06 | True | Adam | 1e-03 | .94 | True | AdamW | 1e-07 | .94 | True | True | |
| oxfordpets | | Adam | 1e-04 | 1e-06 | True | Adam | 1e-03 | .94 | True | AdamW | 1e-07 | .94 | True | True | |
| oxfordpets | | Adam | 1e-04 | 1e-06 | True | Adam | 1e-03 | .94 | True | AdamW | 1e-07 | .94 | True | True | |
| pathmst | | SGD | 1e-01 | 1e-04 | True | Adam | 1e-06 | .93 | True | Adam | 1e-06 | .93 | True | True | |
| pneumoniamst | | AdamW | 1e-03 | 1e-04 | True | Adam | 1e-06 | .93 | True | Adam | 1e-06 | .93 | True | True | |
| retinamst | | SGD | 1e-01 | 1e-04 | True | Adam | 1e-06 | .93 | True | Adam | 1e-06 | .93 | True | True | |
| stanfordcars | | Adam | 1e-03 | 1e-04 | False | Adam | 1e-03 | .99 | False | Adam | 1e-06 | .93 | True | False | |
| tissuemst | SGD | 1e-02 | 1e-04 | True | AdamW | 1e-03 | .94 | True | AdamW | 1e-03 | .93 | True | False | | |
| waterbirds | Adam | 1e-04 | 1e-06 | False | Adam | 1e-06 | .94 | True | Adam | 1e-03 | .93 | True | True | | |
| Dataset | optim. | 1. rate | v. decay | t. decay | mon. SGD | nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. 1. rate | unc. mon. SGD | unc. nest. optim. | unc. b_norm | Arch. | |
| SVHN | SVHN | AdamW | 1e-04 | 1e-04 | False | Adam | 1e-06 | .95 | True | Adam | 1e-06 | .95 | True | False | |
| | catsdogs | AdamW | 1e-03 | 1e-04 | False | SGD | 1e-03 | .95 | True | SGD | 1e-03 | .95 | True | True | |
| | cifar10 | AdamW | 1e-03 | 1e-05 | False | SGD | 1e-03 | .99 | True | SGD | 1e-03 | .99 | True | True | |

Table A12: Best configurations for SR, divided by architectures.

| Dataset | optim. | l. rate | w. decay | t. decay | mom. | SGD | nest. | n.blocks | d.token | d.main | d.hidden | f | d.dropout | res. | dropout | d.ffn | factor | ffn.dropout | Arch. |
|--|--------|---------|----------|----------|-------|-------|-------|----------|---------|--------|----------|------|-----------|------|---------|-------|--------|-------------|-------------|
| adult | SGD | 1e-03 | 1e-04 | 1e-04 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .25 | 2.00 | 2.00 | .25 | FT |
| compass | Adam | 1e-03 | 1e-06 | 1e-06 | False | False | False | 1 | 320 | 1.10 | .10 | 2.00 | .25 | 2.00 | .25 | 2.33 | 2.33 | .35 | Transformer |
| covertype | Adam | 1e-04 | 1e-05 | 1e-05 | False | False | False | 3 | 384 | 1.45 | .45 | 2.33 | .35 | 2.33 | .35 | 2.33 | 2.33 | .35 | Transformer |
| electricity | Adam | 1e-04 | 1e-05 | 1e-05 | False | False | False | 3 | 384 | 1.45 | .45 | 2.33 | .35 | 2.33 | .35 | 2.33 | 2.33 | .35 | Transformer |
| eye | Adam | 1e-05 | 1e-06 | 1e-06 | False | False | False | 1 | 128 | 1.15 | .15 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | FT |
| heloc | SGD | 1e-04 | 1e-04 | 1e-04 | False | True | True | 4 | 128 | .50 | .50 | 2.00 | .20 | 2.00 | .20 | 2.00 | 2.00 | .15 | Transformer |
| higgs | Adam | 1e-05 | 1e-06 | 1e-06 | False | False | False | 2 | 384 | 1.20 | .20 | 2.00 | .15 | 2.00 | .15 | 2.00 | 2.00 | .15 | Transformer |
| house | Adam | 1e-05 | 1e-06 | 1e-06 | False | False | False | 2 | 320 | 1.25 | .25 | 1.00 | .20 | 1.00 | .20 | 1.00 | 1.00 | .20 | Transformer |
| jamnis | Adam | 1e-05 | 1e-06 | 1e-06 | False | False | False | 2 | 384 | 1.35 | .35 | 1.00 | .15 | 1.00 | .15 | 1.00 | 1.00 | .35 | Transformer |
| kdalpmns97 | Adam | 1e-04 | 1e-06 | 1e-06 | False | False | False | 3 | 320 | 1.45 | .45 | 2.33 | .35 | 2.33 | .35 | 2.33 | 2.33 | .35 | Transformer |
| letter | Adam | 1e-04 | 1e-06 | 1e-06 | False | False | False | 1 | 384 | 1.10 | .10 | 1.00 | .35 | 1.00 | .35 | 1.00 | 1.00 | .35 | Transformer |
| magic | Adam | 1e-04 | 1e-06 | 1e-06 | False | False | False | 1 | 384 | 1.10 | .10 | 1.00 | .35 | 1.00 | .35 | 1.00 | 1.00 | .35 | Transformer |
| online | Adam | 1e-02 | 1e-05 | 1e-05 | False | False | False | 1 | 128 | 1.28 | .45 | 1.67 | .45 | 1.67 | .45 | 1.67 | 1.67 | .45 | Transformer |
| phoneme | Adam | 1e-03 | 1e-04 | 1e-04 | False | False | False | 1 | 192 | 1.28 | .45 | 1.67 | .45 | 1.67 | .45 | 1.67 | 1.67 | .45 | Transformer |
| pol | Adam | 1e-03 | 1e-04 | 1e-04 | False | False | False | 1 | 192 | 1.28 | .45 | 1.67 | .45 | 1.67 | .45 | 1.67 | 1.67 | .45 | Transformer |
| rl | Adam | 1e-02 | 1e-03 | 1e-03 | False | False | False | 3 | 512 | .50 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | Transformer |
| upselling | SGD | 1e-04 | 1e-04 | 1e-04 | False | False | False | 4 | 64 | .25 | .25 | 1.00 | .20 | 1.00 | .20 | 1.00 | 1.00 | .20 | Transformer |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. n.blocks d.token d.main d.hidden f d.dropout res. dropout d.ffn factor ffn.dropout Arch. | | | | | | | | | | | | | | | | | | | |
| aloi | AdamW | 1e-05 | 1e-06 | 1e-06 | False | False | False | 4 | 192 | 320 | 4.00 | 1.10 | .20 | 1.00 | .20 | 1.00 | 1.00 | .20 | TabResnet |
| bank | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .30 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | TabResnet |
| giveme | AdamW | 1e-04 | 1e-05 | 1e-05 | True | True | True | 3 | 64 | 256 | 2.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | TabResnet |
| helena | AdamW | 1e-05 | 1e-06 | 1e-06 | False | False | False | 3 | 192 | 320 | 4.00 | 1.10 | .20 | 1.00 | .20 | 1.00 | 1.00 | .20 | TabResnet |
| belena | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | TabResnet |
| miniboone | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | TabResnet |
| uncredit | AdamW | 1e-04 | 1e-04 | 1e-04 | True | True | True | 3 | 192 | 128 | 2.00 | .25 | .25 | 1.00 | .25 | 1.00 | 1.00 | .25 | TabResnet |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. zero_init_resid Arch. | | | | | | | | | | | | | | | | | | | |
| MNIST | | | | | | | | | | | | | | | | | | | |
| bloodmst | Adam | 1e-04 | 1e-04 | 1e-04 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | Resnet18-50 |
| brastmst | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .30 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | Resnet18-50 |
| chestmst | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .30 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | Resnet18-50 |
| dermmst | AdamW | 1e-04 | 1e-03 | 1e-03 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| food101 | Adam | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| ocvmst | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| organmst | SGD | 1e-01 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| organmst | AdamW | 1e-03 | 1e-03 | 1e-03 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| organmst | Adam | 1e-04 | 1e-05 | 1e-05 | False | False | False | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| organmst | Adam | 1e-02 | 1e-05 | 1e-05 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| oxfordipets | SGD | 1e-02 | 1e-03 | 1e-03 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| oxfordipets | Adam | 1e-02 | 1e-05 | 1e-05 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| pathmst | SGD | 1e-04 | 1e-06 | 1e-06 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| pathmst | Adam | 1e-04 | 1e-06 | 1e-06 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| preimniamst | Adam | 1e-04 | 1e-06 | 1e-06 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| preimniamst | Adam | 1e-05 | 1e-06 | 1e-06 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| startordcars | Adam | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| startordcars | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| tismniamst | Adam | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| tismniamst | AdamW | 1e-03 | 1e-04 | 1e-04 | True | True | True | 4 | 256 | 64 | 3.00 | .40 | .40 | 1.00 | .40 | 1.00 | 1.00 | .40 | Resnet18-50 |
| vaterbirds | Adam | 1e-04 | 1e-04 | 1e-04 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | Resnet18-50 |
| Dataset optim. l. rate w. decay t. decay mom. SGD nest. b_norm Arch. | | | | | | | | | | | | | | | | | | | |
| SVHN | | | | | | | | | | | | | | | | | | | |
| casdogs | SGD | 1e-03 | 1e-06 | 1e-06 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | VGG |
| casdogs | Adam | 1e-03 | 1e-04 | 1e-04 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | VGG |
| chrf10 | Adam | 1e-03 | 1e-04 | 1e-04 | False | False | False | 4 | 128 | .25 | .10 | 1.00 | .30 | 1.00 | .30 | 1.00 | 1.00 | .30 | VGG |

Appendix B. Additional Experimental Results

B.1 Q1: Results by Dataset Type

Figure B1 plots the best two and the worst two baselines mean *RelErr* by data type.

For binary tabular datasets (Figure B1a), SAT+EM+SR and SR are the best two performing methods. The former’s relative error rate ranges from $\approx .508$ at $c = .99$ to $\approx .405$ at $c = .70$, while the latter achieves $\approx .511$ at $c = .99$ and $\approx .393$ at $c = .70$. The worst two methods are DG, with *RelErr* of $\approx .632$ at $c = .99$ and $\approx .559$ at $c = .70$, and REG with *RelErr* of $\approx .547$ at $c = .99$ and $\approx .527$ at $c = .70$.

For multiclass tabular datasets (Figure B1c), the best two methods are ENS+SR and SAT+SR, with a mean relative error rate of $\approx .164$ and $\approx .158$ at $c = .99$ respectively, up to $\approx .094$ and $\approx .096$ at $c = .70$. The worst methods are REG, which reaches a mean relative error rate of $\approx .211$ at $c = .99$ and of $\approx .203$ at $c = .70$, and SELNET+EM+SR, with a relative error rate ranging from $\approx .195$ at $c = .99$ to $\approx .218$ at $c = .70$.

For image datasets, methods based on ensembles, i.e., ENS and ENS+SR, achieve the lowest relative error rate. For binary image datasets (Figure B1b), ENS+SR reaches $\approx .378$ at $c = .99$ and $\approx .228$ at $c = .70$, while ENS ranges from $\approx .386$ at $c = .99$ to $\approx .234$ at $c = .70$. In this setting, the worst baselines are SELE and DG, with a mean relative error rate of $\approx .529$ and $\approx .565$ at $c = .99$ respectively, up to $\approx .564$ and $\approx .582$ at $c = .70$ respectively.

For multiclass image datasets (Figure B1d), ENS+SR passes from a mean relative error rate of $\approx .189$ at $c = .99$ to $\approx .126$ at $c = .70$, while ENS achieves $\approx .191$ at $c = .99$ up to $\approx .151$ at $c = .70$. The worst methods here are REG and SELE. The former’s relative error rate ranges from $\approx .276$ at $c = .99$ to $\approx .279$ at $c = .70$, while the latter achieves $\approx .272$ at $c = .99$ and $\approx .238$ at $c = .70$.

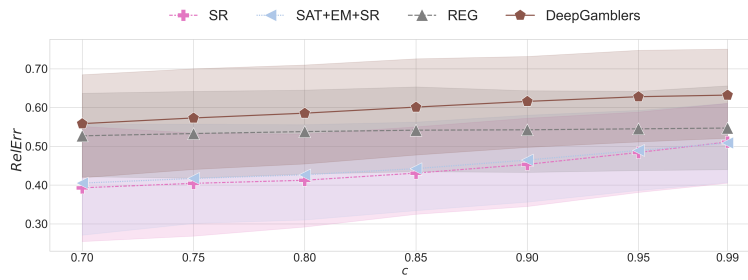
Then, we perform the Nemenyi post hoc test to check for statistically significant differences. Figures B2 and B3 provide CD plots when considering tabular and image data respectively at $c = .99$, $c = .90$, $c = .80$ and $c = .70$. As for aggregated results, all the best performing methods are not distinguishable in a statistically significant sense.

B.2 Q5: Additional Results

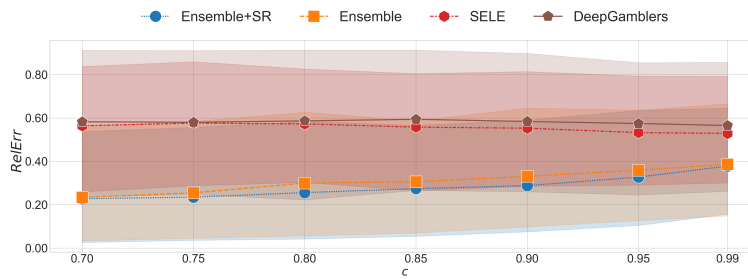
Figure B4 provides the detailed results for the out-of-distribution test sets.

Moreover, we provide additional results w.r.t. distribution shifts. We perform the same experiment as for Q5, but now considering datasets in the `OpenOOD` benchmark (Yang et al., 2022), which is specific for out-of-distribution detection, rather than randomly generated pictures. For `cifar10` we use as test set a random sample from `cifar100`, for `MNIST` a random sample from `FashionMNIST` and for `SVHN` a random sample from `cifar10`. Figure B5 reports the results and the overall mean over the 3 datasets for the 16 baselines considered.

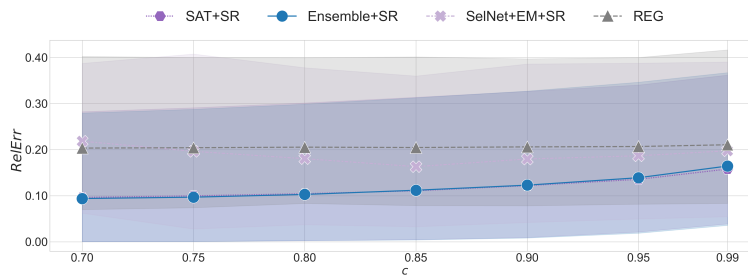
Similarly to the experiments in Section 5.2, we observe that under this milder data shift, no baseline drops all the instances at $c = .99$. We can also see that for lower target coverages, we have higher rejection rates, as expected. Moreover, there is a clear worst-performing method, namely REG.



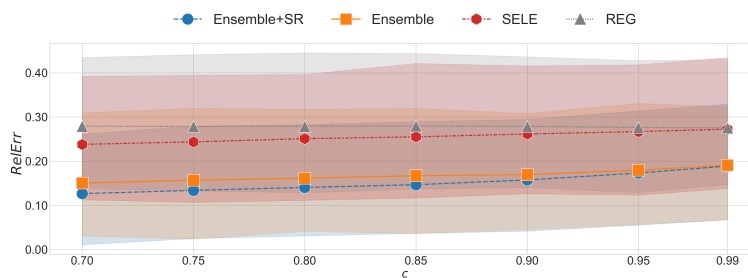
(a) Binary Tabular Data.



(b) Binary Image Data.



(c) Multiclass Tabular Data.



(d) Multiclass Image Data.

Figure B1: Q1: $RelErr$ as a function of target coverage c for the two best and worst approaches on (a) binary tabular data, (b) binary image data, (c) multiclass tabular data and (d) multiclass image data. For readability, only the two best and two worst approaches are shown in each subplot.

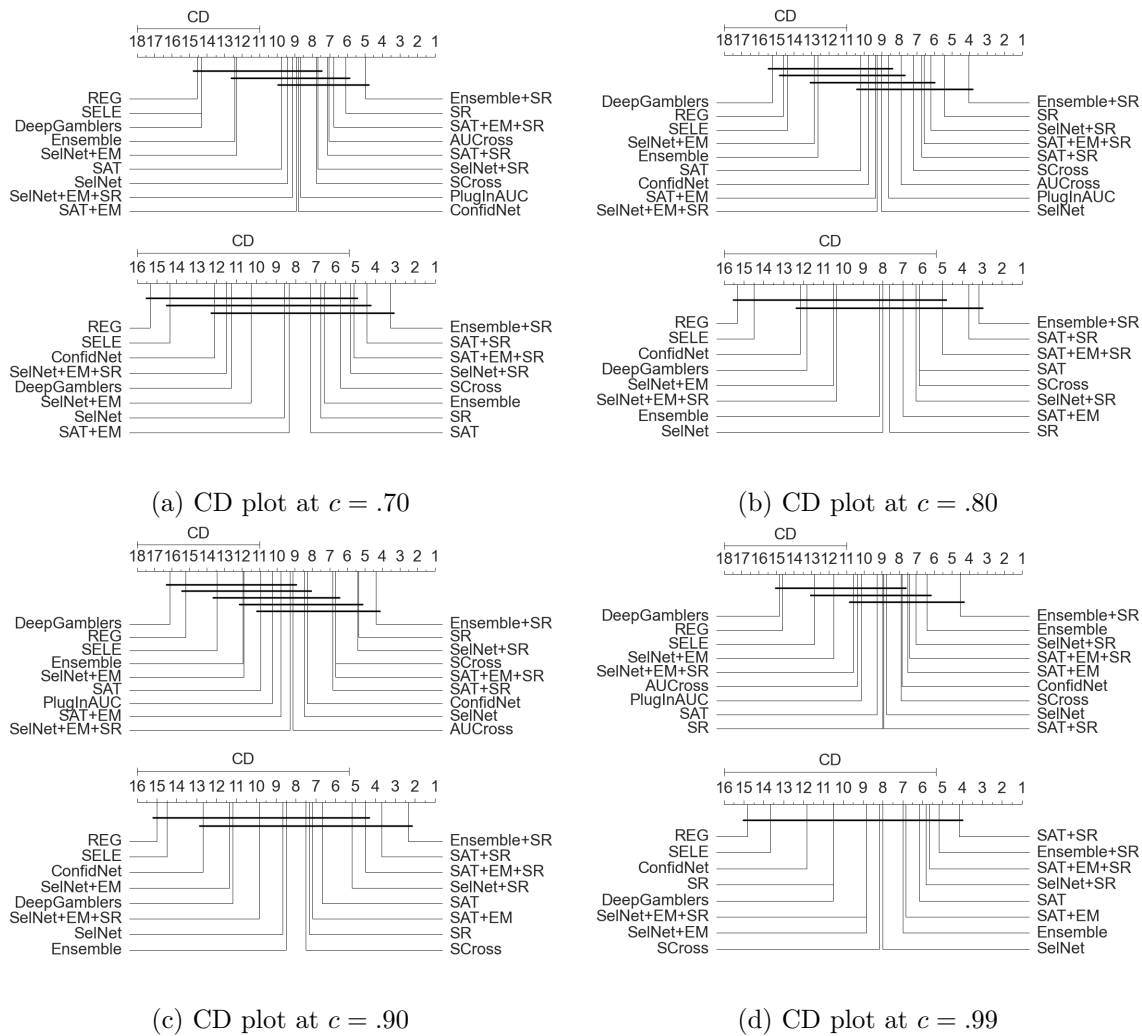


Figure B2: Q1: CD plots of relative error rate $RelErr$ for different target coverages on tabular datasets. Top plots for binary datasets. Bottom plots for multiclass datasets.

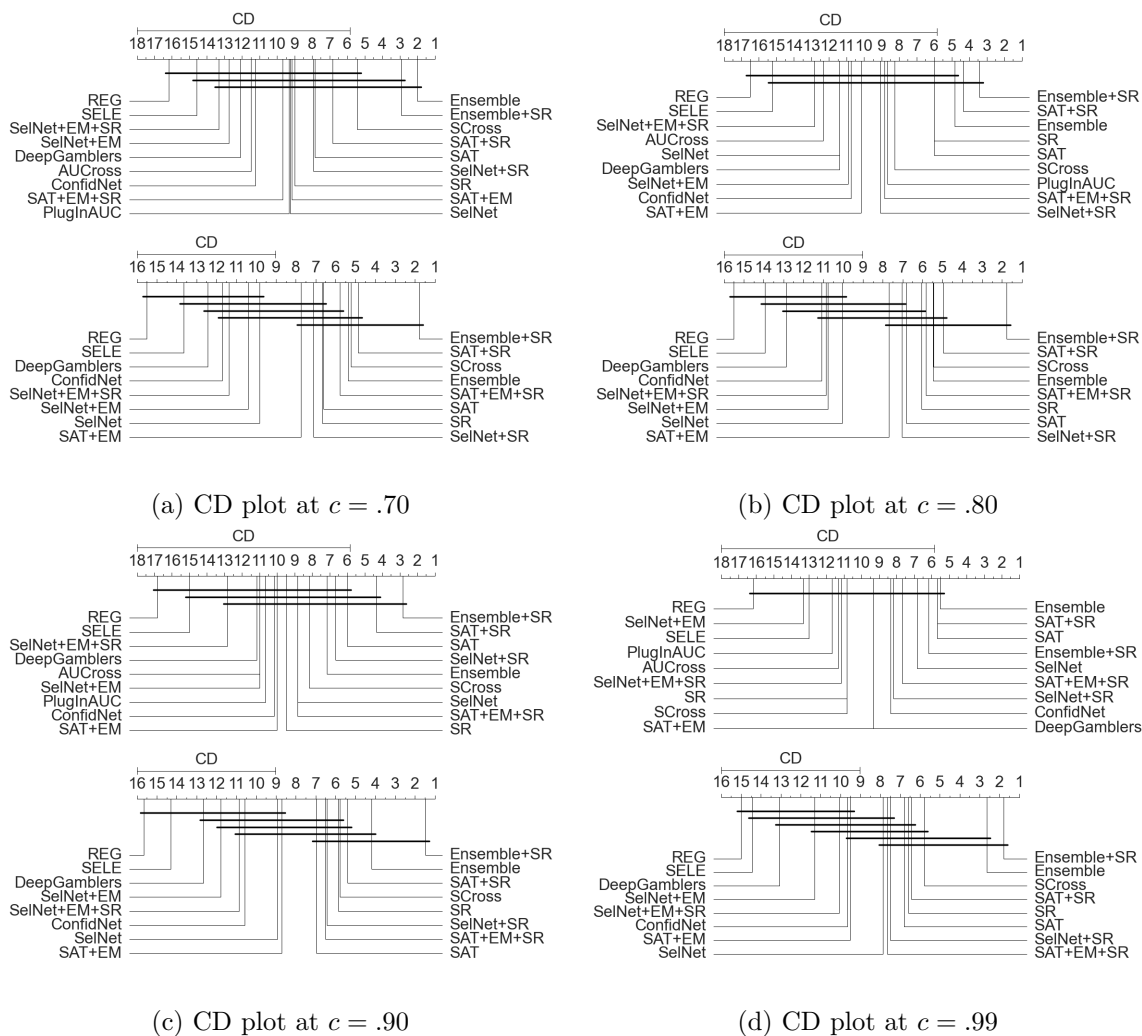


Figure B3: Q1: CD plots of relative error rate $RelErr$ for different target coverages on image datasets. Top plots for binary datasets. Bottom plots for multiclass datasets.

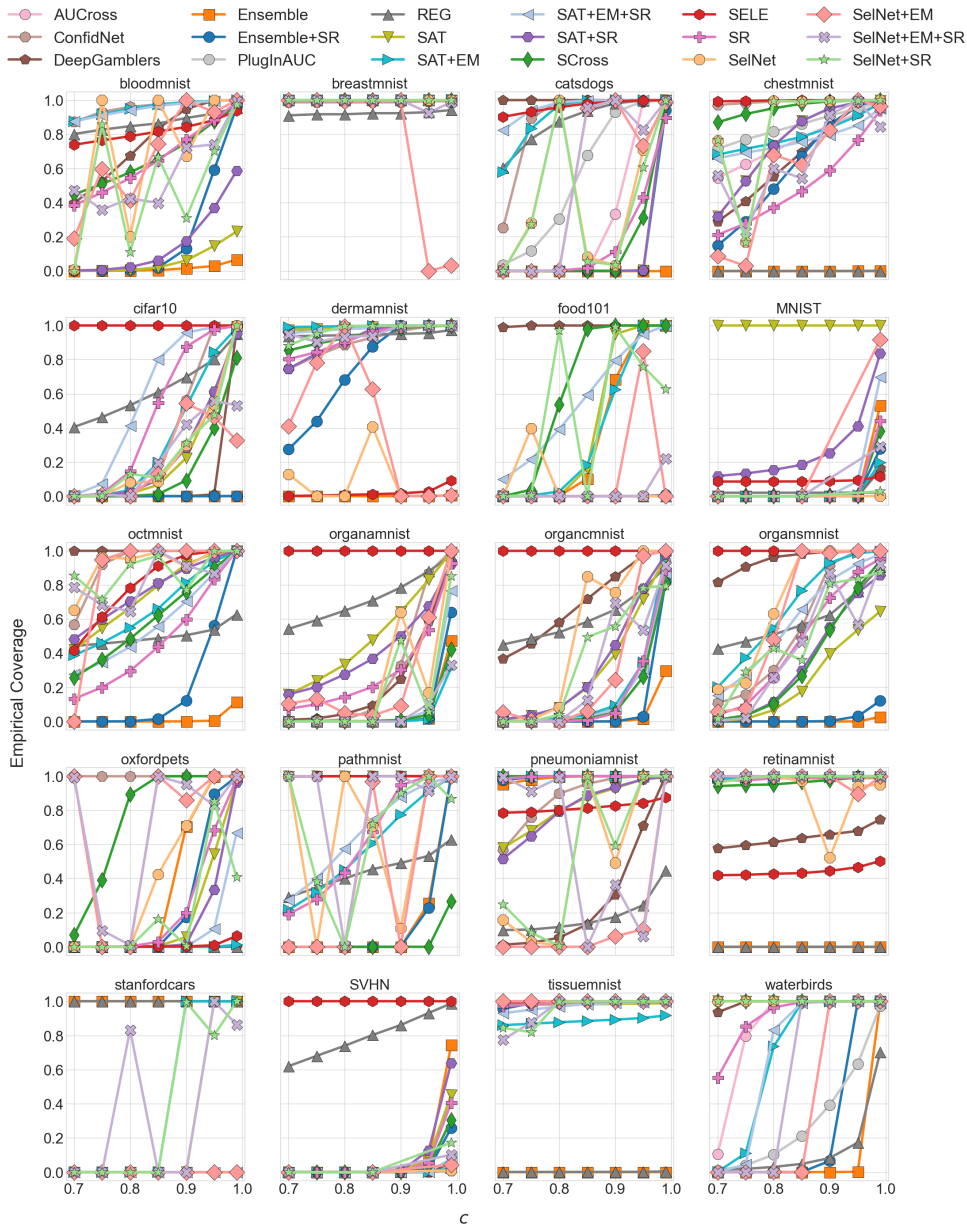


Figure B4: Q5: Empirical coverage $\hat{\phi}$ for out-of-distribution test sets on 20 image datasets for different target coverages c .

For `cifar10`, the method dropping more instances is ENS+SR, reaching an actual coverage of $\approx 5.7\%$ at $c = .70$. The runner-up is ENS, accepting only $\approx 6\%$ of instances. All the remaining baselines have an empirical coverage above 8% at $c = .70$.

For `MNIST` and `SVHN` we observe similar patterns: eleven out of sixteen baselines reach a coverage below 1% at $c = .70$. We also highlight that SELNET reaches zero coverage for $c = .75$ and $.70$ on the `SVHN` dataset.

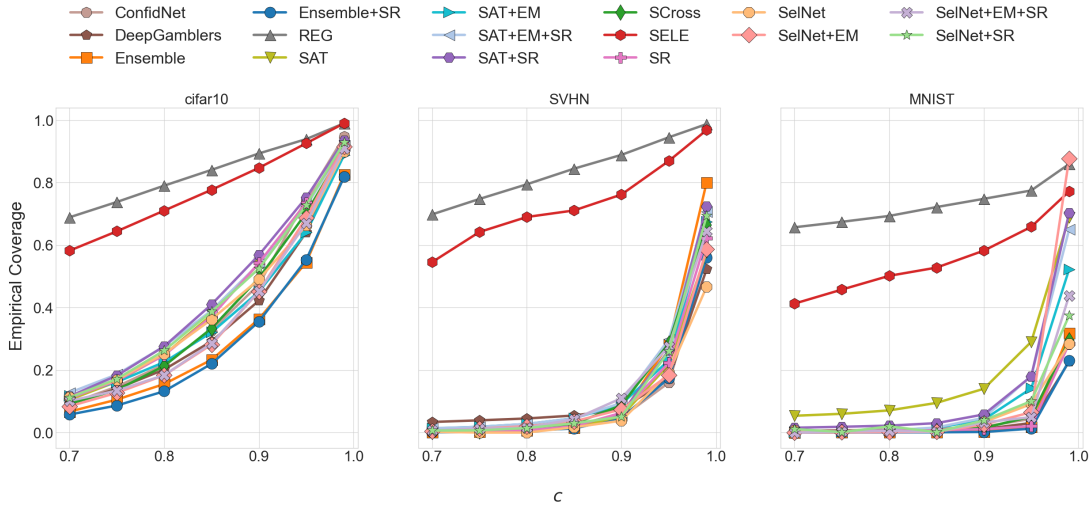


Figure B5: Q5: Empirical coverage $\hat{\phi}$ for out-of-distribution test sets on 3 image datasets for different target coverages c .

To conclude, the experiments show the difficulty for current SC methods to properly perform rejection under distribution shifts. For test data close to the (learned) decision boundary, the baselines correctly reject the instances, since all the methods are built to perform ambiguity rejection. For test data far from the decision boundary, the selection function become confident that the shifted instances are very likely to belong to a certain training class, ending up not rejecting the instances. We think that a potential way to mitigate these problems consists of mixing ambiguity rejection with novelty rejection methods, highlighting the need for further research in this direction.

B.3 Dataset-level Results

We provide detailed results for all the datasets in Tables B1-B44. Each table reports *mean ± std* over the 100 bootstrap samples of \widehat{Err} and empirical coverage $\hat{\phi}$, for every baseline and every target coverage. For binary datasets, we also include the *MinCoeff* metric introduced in the main text. For error rate, we highlight the baseline column with the lowest average (and standard error in case of ties) in bold case. For *Coverage* and *MinCoeff*, we highlight in bold the values closest to target coverage c and 1, respectively.

Table B1: Results for adult: *mean ± std* for \widehat{Err} , empirical coverage $\hat{\phi}$, and *MinCoeff*.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugInAUC |
|-----------------|-----|-------------------|-------------------|------------|-------------------|-------------------|-------------------|------------|-------------------|-------------------|--------------|-------------------|-------------------|-------------------|-------------------|------------|------------|-------------------|-------------------|
| \widehat{Err} | .99 | 153 ± .003 | 136 ± .004 | 125 ± .004 | 123 ± .003 | .160 ± .003 | 123 ± .003 | 132 ± .004 | 132 ± .004 | 132 ± .004 | 158 ± .003 | 130 ± .003 | 125 ± .003 | 138 ± .004 | 133 ± .003 | 130 ± .004 | 133 ± .003 | 130 ± .004 | 131 ± .004 |
| | .95 | 154 ± .003 | 126 ± .004 | 115 ± .004 | 108 ± .003 | .137 ± .004 | 101 ± .003 | 118 ± .004 | 116 ± .004 | 114 ± .004 | 135 ± .003 | 130 ± .003 | 103 ± .003 | 130 ± .004 | 128 ± .004 | 123 ± .003 | 116 ± .003 | 125 ± .004 | 133 ± .004 |
| | .90 | 142 ± .003 | 113 ± .004 | 102 ± .004 | 098 ± .003 | .117 ± .003 | 080 ± .003 | 097 ± .003 | 101 ± .003 | 096 ± .003 | 120 ± .003 | 130 ± .003 | 080 ± .003 | 121 ± .004 | 119 ± .004 | 119 ± .003 | 101 ± .003 | 114 ± .003 | 134 ± .004 |
| | .85 | 119 ± .003 | 097 ± .003 | 088 ± .003 | 080 ± .003 | .106 ± .003 | 061 ± .003 | 080 ± .003 | 083 ± .003 | 078 ± .003 | 106 ± .003 | 128 ± .003 | 061 ± .003 | 106 ± .004 | 101 ± .004 | 113 ± .003 | 086 ± .003 | 103 ± .003 | 131 ± .004 |
| | .80 | 097 ± .003 | 079 ± .003 | 071 ± .003 | 067 ± .003 | 082 ± .003 | 046 ± .002 | 067 ± .003 | 071 ± .003 | 067 ± .003 | 082 ± .003 | 126 ± .003 | 048 ± .003 | 091 ± .003 | 092 ± .004 | 107 ± .004 | 075 ± .003 | 089 ± .003 | 134 ± .004 |
| | .75 | 077 ± .003 | 062 ± .003 | 058 ± .003 | 056 ± .003 | 076 ± .003 | 030 ± .002 | 056 ± .003 | 057 ± .003 | 055 ± .003 | 064 ± .003 | 123 ± .003 | 035 ± .002 | 074 ± .003 | 082 ± .004 | 101 ± .004 | 064 ± .003 | 073 ± .003 | 131 ± .004 |
| | .70 | 063 ± .003 | 055 ± .003 | 046 ± .003 | 045 ± .003 | 066 ± .003 | 020 ± .002 | 044 ± .003 | 044 ± .003 | 045 ± .003 | 065 ± .003 | 118 ± .003 | 024 ± .002 | 056 ± .003 | 076 ± .003 | 099 ± .004 | 055 ± .003 | 057 ± .003 | 128 ± .004 |
| $\hat{\phi}$ | .99 | 983 ± .001 | 994 ± .001 | 967 ± .002 | 924 ± .003 | 998 ± .000 | 980 ± .001 | 985 ± .001 | 991 ± .001 | 991 ± .001 | 994 ± .001 | 995 ± .001 | 986 ± .001 | 980 ± .001 | 984 ± .001 | 983 ± .001 | 992 ± .001 | 991 ± .001 | 990 ± .001 |
| | .95 | 927 ± .003 | 958 ± .002 | 925 ± .002 | 913 ± .003 | 951 ± .002 | 912 ± .003 | 943 ± .002 | 945 ± .002 | 939 ± .002 | 942 ± .002 | 982 ± .001 | 917 ± .002 | 938 ± .003 | 955 ± .002 | 923 ± .003 | 955 ± .002 | 948 ± .002 | 949 ± .002 |
| | .90 | 855 ± .004 | 914 ± .003 | 875 ± .003 | 886 ± .003 | 915 ± .003 | 842 ± .003 | 883 ± .003 | 891 ± .003 | 889 ± .003 | 924 ± .002 | 969 ± .002 | 840 ± .003 | 888 ± .003 | 906 ± .003 | 878 ± .004 | 907 ± .003 | 896 ± .003 | 896 ± .003 |
| | .85 | 793 ± .004 | 855 ± .003 | 825 ± .004 | 818 ± .003 | 854 ± .003 | 771 ± .004 | 824 ± .004 | 836 ± .003 | 820 ± .003 | 825 ± .003 | 960 ± .002 | 767 ± .004 | 839 ± .004 | 813 ± .004 | 859 ± .003 | 844 ± .003 | 844 ± .003 | 836 ± .004 |
| | .80 | 742 ± .004 | 798 ± .004 | 772 ± .004 | 775 ± .004 | 736 ± .004 | 688 ± .004 | 771 ± .004 | 789 ± .004 | 772 ± .004 | 795 ± .004 | 948 ± .002 | 701 ± .004 | 730 ± .004 | 741 ± .004 | 788 ± .004 | 815 ± .004 | 730 ± .004 | 789 ± .004 |
| | .75 | 686 ± .004 | 732 ± .004 | 722 ± .004 | 725 ± .004 | 731 ± .004 | 616 ± .005 | 723 ± .004 | 734 ± .004 | 722 ± .004 | 675 ± .004 | 931 ± .002 | 632 ± .005 | 740 ± .004 | 699 ± .005 | 723 ± .004 | 768 ± .004 | 738 ± .005 | 713 ± .005 |
| | .70 | 632 ± .005 | 679 ± .004 | 671 ± .004 | 672 ± .004 | 683 ± .004 | 563 ± .005 | 669 ± .004 | 679 ± .004 | 674 ± .004 | 677 ± .004 | 909 ± .002 | 579 ± .005 | 687 ± .005 | 670 ± .005 | 686 ± .005 | 728 ± .004 | 687 ± .005 | 659 ± .004 |
| <i>MinCoeff</i> | .99 | 948 ± .018 | 994 ± .019 | 949 ± .020 | 835 ± .017 | 999 ± .019 | 980 ± .019 | 983 ± .019 | 993 ± .018 | 992 ± .019 | 995 ± .019 | 999 ± .019 | 988 ± .019 | 999 ± .019 | 997 ± .019 | 994 ± .019 | 988 ± .019 | 995 ± .019 | 1007 ± .019 |
| | .95 | 788 ± .018 | 946 ± .019 | 872 ± .019 | 897 ± .019 | 941 ± .020 | 907 ± .019 | 933 ± .019 | 942 ± .019 | 934 ± .018 | 929 ± .019 | 977 ± .019 | 921 ± .018 | 988 ± .020 | 981 ± .019 | 960 ± .020 | 938 ± .019 | 988 ± .019 | 1053 ± .020 |
| | .90 | 612 ± .017 | 878 ± .019 | 770 ± .019 | 872 ± .019 | 902 ± .019 | 810 ± .019 | 879 ± .019 | 886 ± .019 | 871 ± .019 | 919 ± .020 | 964 ± .019 | 827 ± .019 | 969 ± .020 | 948 ± .019 | 944 ± .020 | 871 ± .019 | 973 ± .019 | 1068 ± .021 |
| | .85 | 490 ± .016 | 778 ± .017 | 690 ± .016 | 789 ± .018 | 812 ± .018 | 712 ± .018 | 794 ± .018 | 817 ± .018 | 799 ± .018 | 813 ± .018 | 955 ± .019 | 726 ± .018 | 926 ± .019 | 836 ± .019 | 804 ± .020 | 798 ± .019 | 953 ± .020 | 1104 ± .022 |
| | .80 | 393 ± .014 | 675 ± .016 | 604 ± .015 | 729 ± .018 | 801 ± .020 | 606 ± .016 | 730 ± .018 | 762 ± .018 | 739 ± .018 | 800 ± .019 | 943 ± .019 | 637 ± .017 | 872 ± .019 | 720 ± .018 | 854 ± .020 | 731 ± .018 | 933 ± .020 | 1143 ± .022 |
| | .75 | 329 ± .013 | 569 ± .015 | 533 ± .015 | 650 ± .017 | 652 ± .016 | 533 ± .017 | 659 ± .018 | 695 ± .018 | 666 ± .018 | 691 ± .018 | 927 ± .020 | 559 ± .017 | 739 ± .019 | 633 ± .018 | 739 ± .020 | 655 ± .017 | 907 ± .021 | 1184 ± .024 |
| | .70 | 261 ± .013 | 460 ± .015 | 479 ± .016 | 550 ± .016 | 537 ± .015 | 511 ± .018 | 582 ± .017 | 608 ± .016 | 577 ± .017 | 541 ± .015 | 903 ± .019 | 504 ± .017 | 674 ± .018 | 557 ± .017 | 771 ± .020 | 597 ± .017 | 872 ± .021 | 1222 ± .025 |

DEEP NEURAL NETWORK BENCHMARKS FOR SELECTIVE CLASSIFICATION

Table B2: Results for aloi: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | |
|-----------------|----|---------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|---------------|----------------------|----------------------|----------------------|----------------------|----------------------|---------------|
| \widehat{Err} | 99 | 0.035 ± 0.001 | 0.034 ± 0.001 | 0.029 ± 0.001 | 0.039 ± 0.001 | 0.031 ± 0.001 | 0.036 ± 0.001 | 0.033 ± 0.001 | 0.027 ± 0.001 | 0.037 ± 0.001 | 0.030 ± 0.001 | 0.034 ± 0.001 | 0.032 ± 0.001 | 0.035 ± 0.001 | 0.033 ± 0.002 | 0.033 ± 0.002 | 0.031 ± 0.001 | |
| | 95 | 0.026 ± 0.001 | 0.022 ± 0.001 | 0.017 ± 0.001 | 0.032 ± 0.001 | 0.047 ± 0.001 | 0.018 ± 0.001 | 0.017 ± 0.001 | 0.029 ± 0.001 | 0.035 ± 0.001 | 0.035 ± 0.001 | 0.022 ± 0.001 | 0.016 ± 0.001 | 0.022 ± 0.001 | 0.036 ± 0.002 | 0.033 ± 0.002 | 0.016 ± 0.001 | |
| | 90 | 0.022 ± 0.001 | 0.013 ± 0.001 | 0.009 ± 0.001 | 0.026 ± 0.001 | 0.044 ± 0.002 | 0.007 ± 0.001 | 0.007 ± 0.001 | 0.006 ± 0.001 | 0.018 ± 0.001 | 0.024 ± 0.001 | 0.014 ± 0.001 | 0.005 ± 0.001 | 0.030 ± 0.001 | 0.049 ± 0.002 | 0.062 ± 0.002 | 0.006 ± 0.001 | |
| | 85 | 0.019 ± 0.001 | 0.008 ± 0.001 | 0.005 ± 0.000 | 0.023 ± 0.001 | 0.018 ± 0.001 | 0.002 ± 0.000 | 0.002 ± 0.000 | 0.004 ± 0.000 | 0.008 ± 0.001 | 0.012 ± 0.001 | 0.008 ± 0.001 | 0.002 ± 0.000 | 0.030 ± 0.001 | 0.042 ± 0.002 | 0.060 ± 0.002 | 0.002 ± 0.000 | |
| | 80 | 0.016 ± 0.001 | 0.006 ± 0.001 | 0.003 ± 0.000 | 0.022 ± 0.001 | 0.035 ± 0.001 | 0.001 ± 0.000 | 0.001 ± 0.000 | 0.002 ± 0.000 | 0.009 ± 0.001 | 0.011 ± 0.001 | 0.006 ± 0.001 | 0.001 ± 0.000 | 0.030 ± 0.001 | 0.036 ± 0.001 | 0.059 ± 0.002 | 0.001 ± 0.000 | |
| | 75 | 0.014 ± 0.001 | 0.004 ± 0.000 | 0.002 ± 0.000 | 0.019 ± 0.001 | 0.032 ± 0.001 | 0.001 ± 0.000 | 0.001 ± 0.000 | 0.001 ± 0.000 | 0.006 ± 0.001 | 0.011 ± 0.001 | 0.002 ± 0.000 | 0.000 ± 0.000 | 0.029 ± 0.001 | 0.031 ± 0.001 | 0.056 ± 0.002 | 0.001 ± 0.000 | |
| | 70 | 0.012 ± 0.001 | 0.003 ± 0.000 | 0.002 ± 0.000 | 0.019 ± 0.001 | 0.022 ± 0.001 | 0.000 ± 0.000 | 0.000 ± 0.000 | 0.001 ± 0.000 | 0.005 ± 0.001 | 0.012 ± 0.001 | 0.001 ± 0.000 | 0.000 ± 0.000 | 0.028 ± 0.001 | 0.027 ± 0.001 | 0.054 ± 0.002 | 0.000 ± 0.000 | |
| | 99 | 0.991 ± 0.001 | 0.989 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 | 0.989 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 | 0.990 ± 0.001 | 0.991 ± 0.001 | 0.992 ± 0.001 | 0.991 ± 0.001 | 0.990 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 | 0.994 ± 0.001 |
| | 95 | 0.951 ± 0.001 | 0.950 ± 0.001 | 0.954 ± 0.001 | 0.949 ± 0.001 | 0.951 ± 0.001 | 0.951 ± 0.001 | 0.951 ± 0.001 | 0.951 ± 0.001 | 0.956 ± 0.001 | 0.954 ± 0.001 | 0.952 ± 0.001 | 0.950 ± 0.001 | 0.952 ± 0.001 | 0.952 ± 0.001 | 0.948 ± 0.001 | 0.959 ± 0.001 | |
| | 90 | 0.904 ± 0.002 | 0.900 ± 0.002 | 0.902 ± 0.002 | 0.901 ± 0.002 | 0.901 ± 0.002 | 0.902 ± 0.002 | 0.903 ± 0.002 | 0.903 ± 0.002 | 0.905 ± 0.002 | 0.909 ± 0.002 | 0.907 ± 0.002 | 0.909 ± 0.002 | 0.901 ± 0.002 | 0.901 ± 0.002 | 0.902 ± 0.002 | 0.900 ± 0.002 | 0.916 ± 0.002 |

Table B3: Results for bank: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhgInAUC | |
|-----------------|----|----------------------|----------------------|----------------------|----------------------|---------------|----------------------|----------------------|----------------------|---------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| \widehat{Err} | 99 | 0.022 ± 0.003 | 0.009 ± 0.002 | 0.007 ± 0.002 | 0.022 ± 0.002 | 0.022 ± 0.002 | 0.027 ± 0.002 | 0.087 ± 0.002 | 0.087 ± 0.002 | 0.090 ± 0.002 | 0.092 ± 0.002 | 0.092 ± 0.003 | 0.090 ± 0.002 | 0.084 ± 0.002 | 0.091 ± 0.002 | 0.103 ± 0.003 | 0.090 ± 0.002 | 0.092 ± 0.002 | 0.095 ± 0.003 | |
| | 95 | 0.010 ± 0.003 | 0.073 ± 0.002 | 0.073 ± 0.002 | 0.076 ± 0.002 | 0.075 ± 0.002 | 0.072 ± 0.002 | 0.068 ± 0.002 | 0.072 ± 0.002 | 0.073 ± 0.002 | 0.073 ± 0.002 | 0.072 ± 0.002 | 0.073 ± 0.002 | 0.072 ± 0.002 | 0.071 ± 0.002 | 0.081 ± 0.002 | 0.103 ± 0.003 | 0.089 ± 0.003 | 0.093 ± 0.003 | |
| | 90 | 0.003 ± 0.002 | 0.054 ± 0.002 | 0.055 ± 0.002 | 0.056 ± 0.002 | 0.056 ± 0.002 | 0.054 ± 0.002 | 0.051 ± 0.002 | 0.053 ± 0.002 | 0.053 ± 0.002 | 0.053 ± 0.002 | 0.053 ± 0.002 | 0.053 ± 0.002 | 0.054 ± 0.002 | 0.055 ± 0.002 | 0.063 ± 0.002 | 0.105 ± 0.003 | 0.055 ± 0.002 | 0.084 ± 0.003 | |
| | 85 | 0.014 ± 0.002 | 0.038 ± 0.002 | 0.039 ± 0.002 | 0.042 ± 0.002 | 0.041 ± 0.002 | 0.039 ± 0.002 | 0.036 ± 0.002 | 0.037 ± 0.002 | 0.038 ± 0.002 | 0.042 ± 0.002 | 0.042 ± 0.002 | 0.042 ± 0.002 | 0.043 ± 0.002 | 0.043 ± 0.002 | 0.050 ± 0.002 | 0.103 ± 0.003 | 0.036 ± 0.002 | 0.078 ± 0.003 | 0.087 ± 0.003 |
| | 80 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.027 ± 0.002 | 0.028 ± 0.002 | 0.027 ± 0.002 | 0.029 ± 0.002 | 0.025 ± 0.001 | 0.027 ± 0.002 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.028 ± 0.002 | 0.034 ± 0.002 | 0.101 ± 0.003 | 0.025 ± 0.002 | 0.069 ± 0.003 | 0.081 ± 0.003 |
| | 75 | 0.019 ± 0.001 | 0.016 ± 0.001 | 0.018 ± 0.001 | 0.019 ± 0.001 | 0.018 ± 0.001 | 0.019 ± 0.001 | 0.017 ± 0.001 | 0.017 ± 0.001 | 0.019 ± 0.001 | 0.022 ± 0.001 | 0.022 ± 0.001 | 0.023 ± 0.002 | 0.018 ± 0.001 | 0.019 ± 0.002 | 0.021 ± 0.002 | 0.099 ± 0.003 | 0.017 ± 0.001 | 0.052 ± 0.003 | 0.071 ± 0.003 |
| | 70 | 0.013 ± 0.001 | 0.013 ± 0.001 | 0.013 ± 0.001 | 0.013 ± 0.001 | 0.014 ± 0.001 | 0.012 ± 0.001 | 0.012 ± 0.001 | 0.011 ± 0.001 | 0.013 ± 0.001 | 0.015 ± 0.001 | 0.015 ± 0.001 | 0.023 ± 0.002 | 0.012 ± 0.001 | 0.012 ± 0.001 | 0.015 ± 0.001 | 0.099 ± 0.003 | 0.012 ± 0.001 | 0.037 ± 0.002 | 0.056 ± 0.003 |
| | 99 | 0.990 ± 0.001 | 0.990 ± 0.001 | 0.988 ± 0.001 | 0.992 ± 0.001 | 0.988 ± 0.001 | 0.991 ± 0.001 | 0.987 ± 0.001 | 0.986 ± 0.001 | 0.988 ± 0.001 | 0.991 ± 0.001 | 0.989 ± 0.001 | 0.989 ± 0.001 | 0.986 ± 0.001 | 0.990 ± 0.001 | 0.990 ± 0.001 | 0.991 ± 0.001 | 0.990 ± 0.001 | 0.991 ± 0.001 | 0.991 ± 0.001 |
| | 95 | 0.952 ± 0.002 | 0.947 ± 0.002 | 0.950 ± 0.002 | 0.949 ± 0.002 | 0.947 ± 0.002 | 0.942 ± 0.002 | 0.941 ± 0.002 | 0.946 ± 0.002 | 0.946 ± 0.002 | 0.946 ± 0.002 | 0.946 ± 0.002 | 0.946 ± 0.002 | 0.946 ± 0.002 | 0.951 ± 0.002 | 0.949 ± 0.002 | 0.949 ± 0.002 | 0.950 ± 0.002 | 0.951 ± 0.002 | 0.948 ± 0.002 |
| | 90 | 0.908 ± 0.003 | 0.905 ± 0.003 | 0.897 ± 0.003 | 0.899 ± 0.003 | 0.897 ± 0.003 | 0.896 ± 0.003 | 0.892 ± 0.003 | 0.892 ± 0.003 | 0.890 ± 0.003 | 0.892 ± 0.003 | 0.893 ± 0.003 | 0.893 ± 0.003 | 0.894 ± 0.003 | 0.894 ± 0.003 | 0.896 ± 0.003 | 0.901 ± 0.003 | 0.908 ± 0.003 | 0.900 ± 0.003 | 0.899 ± 0.003 |

Table B4: Results for bloodmnist: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | | | |
|-----------------|----|----------------------|----------------------|---------------|----------------------|---------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|---------------|----------------------|----------------------|---------------|---------------|----------------------|----------------------|
| \widehat{Err} | 99 | 0.066 ± 0.004 | 0.033 ± 0.003 | 0.039 ± 0.003 | 0.036 ± 0.003 | 0.045 ± 0.003 | 0.043 ± 0.003 | 0.036 ± 0.003 | 0.039 ± 0.003 | 0.038 ± 0.003 | 0.043 ± 0.003 | 0.037 ± 0.003 | 0.036 ± 0.003 | 0.037 ± 0.003 | 0.036 ± 0.003 | 0.064 ± 0.004 | 0.064 ± 0.004 | 0.041 ± 0.003 | | |
| | 95 | 0.056 ± 0.004 | 0.028 ± 0.003 | 0.026 ± 0.003 | 0.025 ± 0.003 | 0.034 ± 0.003 | 0.024 ± 0.003 | 0.022 ± 0.003 | 0.026 ± 0.003 | 0.021 ± 0.002 | 0.030 ± 0.003 | 0.022 ± 0.002 | 0.019 ± 0.002 | 0.035 ± 0.003 | 0.063 ± 0.004 | 0.065 ± 0.004 | 0.029 ± 0.003 | | | |
| | 90 | 0.043 ± 0.003 | 0.016 ± 0.002 | 0.014 ± 0.002 | 0.015 ± 0.002 | 0.021 ± 0.003 | 0.015 ± 0.002 | 0.014 ± 0.002 | 0.013 ± 0.002 | 0.013 ± 0.002 | 0.018 ± 0.002 | 0.013 ± 0.002 | 0.012 ± 0.002 | 0.027 ± 0.004 | 0.062 ± 0.004 | 0.065 ± 0.004 | 0.015 ± 0.002 | | | |
| | 85 | 0.037 ± 0.003 | 0.010 ± 0.002 | 0.008 ± 0.002 | 0.012 ± 0.002 | 0.009 ± 0.002 | 0.010 ± 0.002 | 0.008 ± 0.002 | 0.007 ± 0.002 | 0.008 ± 0.002 | 0.008 ± 0.002 | 0.008 ± 0.002 | 0.008 ± 0.002 | 0.008 ± 0.002 | 0.008 ± 0.002 | 0.020 ± 0.004 | 0.063 ± 0.004 | 0.063 ± 0.004 | 0.009 ± 0.002 | |
| | 80 | 0.029 ± 0.003 | 0.005 ± 0.001 | 0.005 ± 0.002 | 0.007 ± 0.002 | 0.011 ± 0.002 | 0.005 ± 0.001 | 0.005 ± 0.001 | 0.005 ± 0.001 | 0.005 ± 0.001 | 0.005 ± 0.001 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.016 ± 0.004 | 0.064 ± 0.004 | 0.065 ± 0.004 | 0.005 ± 0.002 | |
| | 75 | 0.025 ± 0.003 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.006 ± 0.002 | 0.006 ± 0.002 | 0.004 ± 0.001 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.004 ± 0.001 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.002 ± 0.001 | 0.011 ± 0.004 | 0.065 ± 0.004 | 0.066 ± 0.004 | 0.066 ± 0.004 | 0.003 ± 0.001 | | |
| | 70 | 0.018 ± 0.003 | 0.004 ± 0.001 | 0.002 ± 0.001 | 0.003 ± 0.001 | 0.005 ± 0.001 | 0.004 ± 0.001 | 0.003 ± 0.001 | 0.003 ± 0.001 | 0.003 ± 0.001 | 0.004 ± 0.001 | 0.003 ± 0.001 | 0.004 ± 0.001 | 0.003 ± 0.001 | 0.003 ± 0.001 | 0.003 ± 0.001 | 0.065 ± 0.004 | 0.065 ± 0.004 | 0.002 ± 0.001 | |
| | 99 | 0.990 ± 0.002 | 0.984 ± 0.002 | 0.991 ± 0.002 | 0.990 ± 0.002 | 0.989 ± 0.002 | 0.992 ± 0.002 | 0.991 ± 0.002 | 0.991 ± 0.002 | 0.991 ± 0.002 | 0.992 ± 0.002 | 0.992 ± 0.002 | 0.992 ± 0.002 | 0.988 ± 0.002 | 0.989 ± 0.002 | 0.989 ± 0.002 | 0.991 ± 0.002 | 0.991 ± 0.002 | 0.993 ± 0.002 | 0.991 ± 0.001 |
| | 95 | 0.959 ± 0.003 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.948 ± 0.004 | 0.947 ± 0.004 | 0.948 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.946 ± 0.004 | 0.949 ± 0.004 | 0.941 ± 0.004 | 0.942 ± 0.004 | 0.947 ± 0.004 | 0.950 ± 0.004 | 0.955 ± 0.004 |
| | 90 | 0.909 ± 0.005 | 0.886 ± 0.005 | 0.894 ± 0.005 | 0.900 ± 0.005 | 0.898 ± 0.004 | 0.901 ± 0.006 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.891 ± 0.005 | 0.901 ± 0.006 | 0.903 ± 0.006 | 0.901 ± 0.005 | 0.907 ± 0.005 |

Table B5: Results for breastm18: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhgInAUC |
|-----------------|----|----------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|----------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| \widehat{Err} | 99 | 0.137 ± 0.029 | 0.161 ± 0.031 | 0.159 ± 0.030 | 0.143 ± 0.031 | 0.177 ± 0.033 | 0.176 ± 0.034 | 0.158 ± 0.031 | 0.145 ± 0.030 | 0.144 ± 0.031 | 0.163 ± 0.033 | 0.163 ± 0.032 | 0.151 ± 0.032 | 0.189 ± 0.035 | 0.189 ± 0.033 | 0.205 ± 0.035 | 0.149 ± 0.031 | 0.147 ± 0.031 | 0.188 ± 0.034 |
| | 95 | 0.145 ± 0.031 | 0.160 ± 0.031 | 0.158 ± 0.031 | 0.154 ± 0.033 | 0.192 ± 0.034 | 0.164 ± 0.033 | 0.154 ± 0.032 | 0.129 ± 0.027 | 0.119 ± 0.029 | 0.170 ± 0.034 | 0.145 ± 0.031 | 0.122 ± 0.031 | 0.173 ± 0.037 | 0.186 ± 0.034 | 0.206 ± 0.036 | 0.138 ± 0.030 | 0.146 ± 0.032 | 0.184 ± 0.034 |
| | 90 | 0.150 ± 0.031 | 0.145 ± 0.031 | 0.135 ± 0.031 | 0.157 ± 0.035 | 0.155 ± 0.033 | 0.161 ± 0.034 | | | | | | | | | | | | |

DEEP NEURAL NETWORK BENCHMARKS FOR SELECTIVE CLASSIFICATION

Table B11: Results for dermannist: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross |
|-----------------|-----|------------|-------------------|------------|------------|------------|-------------------|-------------------|-------------------|------------|-------------|-------------------|-------------------|-------------------|------------|------------|-------------------|
| \widehat{Err} | .99 | 269 ± 0.10 | 229 ± 0.09 | 239 ± 0.09 | 243 ± 0.10 | 241 ± 0.10 | 228 ± 0.09 | 236 ± 0.09 | 239 ± 0.11 | 240 ± 0.10 | 226 ± 0.09 | 223 ± 0.09 | 245 ± 0.10 | 270 ± 0.09 | 273 ± 0.09 | 231 ± 0.09 | 231 ± 0.09 |
| | .95 | 250 ± 0.10 | 217 ± 0.09 | 230 ± 0.09 | 227 ± 0.10 | 227 ± 0.10 | 207 ± 0.09 | 214 ± 0.09 | 222 ± 0.09 | 224 ± 0.11 | 213 ± 0.08 | 206 ± 0.09 | 238 ± 0.10 | 263 ± 0.10 | 267 ± 0.10 | 214 ± 0.08 | 214 ± 0.08 |
| | .90 | 215 ± 0.10 | 202 ± 0.09 | 212 ± 0.10 | 207 ± 0.10 | 210 ± 0.10 | 191 ± 0.09 | 194 ± 0.09 | 196 ± 0.09 | 195 ± 0.11 | 206 ± 0.10 | 194 ± 0.09 | 182 ± 0.09 | 218 ± 0.10 | 259 ± 0.10 | 263 ± 0.11 | 193 ± 0.08 |
| | .85 | 194 ± 0.10 | 180 ± 0.09 | 201 ± 0.09 | 198 ± 0.10 | 198 ± 0.10 | 177 ± 0.09 | 173 ± 0.09 | 174 ± 0.09 | 182 ± 0.10 | 180 ± 0.10 | 175 ± 0.09 | 161 ± 0.09 | 200 ± 0.10 | 252 ± 0.10 | 263 ± 0.11 | 169 ± 0.09 |
| | .80 | 180 ± 0.10 | 159 ± 0.09 | 179 ± 0.09 | 192 ± 0.10 | 183 ± 0.11 | 152 ± 0.09 | 153 ± 0.09 | 162 ± 0.09 | 202 ± 0.10 | 199 ± 0.10 | 159 ± 0.10 | 147 ± 0.08 | 179 ± 0.10 | 242 ± 0.11 | 261 ± 0.11 | 146 ± 0.08 |
| | .75 | 163 ± 0.10 | 144 ± 0.09 | 162 ± 0.10 | 164 ± 0.11 | 165 ± 0.10 | 129 ± 0.09 | 129 ± 0.09 | 140 ± 0.09 | 182 ± 0.10 | 185 ± 0.10 | 146 ± 0.10 | 123 ± 0.08 | 161 ± 0.11 | 224 ± 0.11 | 258 ± 0.12 | 128 ± 0.08 |
| | .70 | 143 ± 0.11 | 120 ± 0.09 | 143 ± 0.09 | 146 ± 0.10 | 148 ± 0.10 | 114 ± 0.08 | 110 ± 0.09 | 119 ± 0.09 | 172 ± 0.10 | 165 ± 0.10 | 130 ± 0.10 | 108 ± 0.08 | 147 ± 0.10 | 203 ± 0.10 | 254 ± 0.11 | 109 ± 0.09 |
| $\hat{\phi}$ | .99 | 994 ± 0.02 | 983 ± 0.03 | 993 ± 0.02 | 993 ± 0.02 | 992 ± 0.02 | 991 ± 0.02 | 989 ± 0.02 | 990 ± 0.02 | 990 ± 0.02 | 992 ± 0.02 | 994 ± 0.02 | 994 ± 0.02 | 991 ± 0.02 | 989 ± 0.02 | 993 ± 0.02 | 989 ± 0.02 |
| | .95 | 962 ± 0.04 | 946 ± 0.04 | 965 ± 0.04 | 955 ± 0.05 | 964 ± 0.04 | 947 ± 0.05 | 949 ± 0.05 | 950 ± 0.05 | 954 ± 0.04 | 959 ± 0.04 | 956 ± 0.05 | 960 ± 0.04 | 967 ± 0.04 | 933 ± 0.07 | 957 ± 0.05 | 952 ± 0.04 |
| | .90 | 898 ± 0.08 | 900 ± 0.06 | 906 ± 0.07 | 912 ± 0.06 | 920 ± 0.06 | 905 ± 0.06 | 896 ± 0.07 | 892 ± 0.08 | 898 ± 0.07 | 907 ± 0.07 | 896 ± 0.06 | 910 ± 0.06 | 913 ± 0.07 | 892 ± 0.08 | 905 ± 0.06 | 905 ± 0.06 |
| | .85 | 862 ± 0.08 | 849 ± 0.08 | 868 ± 0.07 | 866 ± 0.08 | 882 ± 0.07 | 865 ± 0.07 | 851 ± 0.08 | 845 ± 0.08 | 852 ± 0.08 | 857 ± 0.07 | 837 ± 0.08 | 852 ± 0.08 | 871 ± 0.08 | 850 ± 0.08 | 870 ± 0.08 | 849 ± 0.07 |
| | .80 | 820 ± 0.09 | 785 ± 0.09 | 818 ± 0.08 | 820 ± 0.09 | 823 ± 0.09 | 797 ± 0.09 | 796 ± 0.09 | 816 ± 0.09 | 800 ± 0.10 | 807 ± 0.10 | 789 ± 0.08 | 809 ± 0.08 | 818 ± 0.09 | 803 ± 0.08 | 818 ± 0.09 | 800 ± 0.08 |
| | .75 | 760 ± 0.09 | 746 ± 0.10 | 757 ± 0.08 | 757 ± 0.10 | 764 ± 0.09 | 739 ± 0.11 | 749 ± 0.09 | 760 ± 0.11 | 760 ± 0.11 | 763 ± 0.10 | 730 ± 0.09 | 750 ± 0.09 | 766 ± 0.10 | 777 ± 0.09 | 777 ± 0.09 | 749 ± 0.09 |
| | .70 | 709 ± 0.10 | 687 ± 0.10 | 705 ± 0.10 | 711 ± 0.11 | 705 ± 0.10 | 702 ± 0.10 | 692 ± 0.10 | 698 ± 0.11 | 720 ± 0.10 | 703 ± 0.10 | 684 ± 0.09 | 702 ± 0.10 | 711 ± 0.11 | 694 ± 0.10 | 736 ± 0.10 | 681 ± 0.10 |

Table B12: Results for electricity: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugInAUC | |
|-----------------|-----|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|-------------------|---------------------|---------------------|-------------------|---------------------|---------------------|-------------------|---------------------|---------------------|-------------------|---------------------|
| \widehat{Err} | .99 | 248 ± 0.04 | 176 ± 0.05 | 183 ± 0.05 | 179 ± 0.04 | 187 ± 0.05 | 174 ± 0.04 | 176 ± 0.05 | 162 ± 0.05 | 179 ± 0.04 | 185 ± 0.05 | 162 ± 0.05 | 161 ± 0.05 | 167 ± 0.04 | 206 ± 0.05 | 200 ± 0.04 | 169 ± 0.05 | 170 ± 0.05 | 174 ± 0.04 | |
| | .95 | 245 ± 0.04 | 164 ± 0.04 | 152 ± 0.04 | 172 ± 0.05 | 174 ± 0.04 | 162 ± 0.04 | 162 ± 0.04 | 149 ± 0.05 | 166 ± 0.04 | 170 ± 0.04 | 159 ± 0.05 | 162 ± 0.04 | 155 ± 0.04 | 198 ± 0.05 | 196 ± 0.04 | 158 ± 0.04 | 161 ± 0.05 | 166 ± 0.04 | |
| | .90 | 228 ± 0.05 | 151 ± 0.04 | 138 ± 0.04 | 161 ± 0.04 | 162 ± 0.05 | 147 ± 0.04 | 145 ± 0.04 | 135 ± 0.04 | 155 ± 0.05 | 160 ± 0.04 | 155 ± 0.05 | 134 ± 0.04 | 139 ± 0.04 | 190 ± 0.05 | 193 ± 0.04 | 144 ± 0.04 | 148 ± 0.05 | 158 ± 0.05 | |
| | .85 | 227 ± 0.05 | 135 ± 0.04 | 125 ± 0.04 | 149 ± 0.05 | 143 ± 0.05 | 135 ± 0.04 | 133 ± 0.04 | 120 ± 0.04 | 139 ± 0.05 | 139 ± 0.04 | 148 ± 0.05 | 120 ± 0.04 | 129 ± 0.04 | 182 ± 0.05 | 190 ± 0.05 | 130 ± 0.04 | 136 ± 0.04 | 143 ± 0.04 | |
| | .80 | 217 ± 0.05 | 121 ± 0.04 | 112 ± 0.04 | 136 ± 0.04 | 138 ± 0.05 | 119 ± 0.04 | 117 ± 0.04 | 107 ± 0.04 | 125 ± 0.04 | 138 ± 0.05 | 140 ± 0.05 | 107 ± 0.04 | 117 ± 0.04 | 178 ± 0.05 | 187 ± 0.05 | 117 ± 0.04 | 121 ± 0.04 | 139 ± 0.05 | |
| | .75 | 203 ± 0.05 | 110 ± 0.04 | 102 ± 0.04 | 125 ± 0.05 | 125 ± 0.05 | 107 ± 0.04 | 101 ± 0.04 | 996 ± 0.04 | 113 ± 0.04 | 125 ± 0.05 | 131 ± 0.04 | 093 ± 0.04 | 109 ± 0.04 | 169 ± 0.05 | 184 ± 0.05 | 104 ± 0.04 | 106 ± 0.04 | 118 ± 0.04 | |
| | .70 | 186 ± 0.05 | 96 ± 0.04 | 98 ± 0.04 | 124 ± 0.04 | 116 ± 0.05 | 97 ± 0.04 | 89 ± 0.04 | 084 ± 0.04 | 112 ± 0.05 | 117 ± 0.05 | 127 ± 0.05 | 085 ± 0.04 | 098 ± 0.04 | 102 ± 0.05 | 181 ± 0.05 | 99 ± 0.04 | 102 ± 0.05 | 105 ± 0.04 | |
| $\hat{\phi}$ | .99 | 990 ± 0.01 | 991 ± 0.01 | 992 ± 0.01 | 992 ± 0.01 | 995 ± 0.01 | 990 ± 0.01 | 989 ± 0.01 | 990 ± 0.01 | 990 ± 0.01 | 992 ± 0.01 | 990 ± 0.01 | 992 ± 0.01 | 988 ± 0.01 | 990 ± 0.01 | 989 ± 0.01 | 992 ± 0.01 | 990 ± 0.01 | 989 ± 0.01 | |
| | .95 | 909 ± 0.03 | 905 ± 0.03 | 953 ± 0.02 | 961 ± 0.02 | 958 ± 0.02 | 950 ± 0.03 | 949 ± 0.03 | 954 ± 0.02 | 953 ± 0.02 | 953 ± 0.03 | 955 ± 0.02 | 951 ± 0.03 | 950 ± 0.03 | 953 ± 0.02 | 948 ± 0.03 | 954 ± 0.02 | 954 ± 0.02 | 951 ± 0.02 | |
| | .90 | 898 ± 0.04 | 901 ± 0.04 | 902 ± 0.04 | 911 ± 0.03 | 907 ± 0.04 | 898 ± 0.03 | 896 ± 0.04 | 905 ± 0.04 | 908 ± 0.03 | 902 ± 0.04 | 909 ± 0.03 | 900 ± 0.04 | 901 ± 0.04 | 908 ± 0.04 | 900 ± 0.03 | 900 ± 0.03 | 900 ± 0.03 | 903 ± 0.04 | |
| | .85 | 849 ± 0.04 | 850 ± 0.04 | 859 ± 0.04 | 859 ± 0.04 | 852 ± 0.04 | 856 ± 0.04 | 851 ± 0.04 | 852 ± 0.04 | 854 ± 0.04 | 849 ± 0.04 | 860 ± 0.04 | 852 ± 0.04 | 853 ± 0.04 | 853 ± 0.05 | 844 ± 0.04 | 863 ± 0.04 | 856 ± 0.04 | 856 ± 0.05 | |
| | .80 | 801 ± 0.04 | 803 ± 0.05 | 817 ± 0.05 | 817 ± 0.04 | 807 ± 0.05 | 807 ± 0.05 | 798 ± 0.05 | 805 ± 0.05 | 806 ± 0.04 | 801 ± 0.05 | 807 ± 0.04 | 803 ± 0.05 | 802 ± 0.05 | 796 ± 0.05 | 795 ± 0.05 | 813 ± 0.05 | 808 ± 0.05 | 802 ± 0.05 | |
| | .75 | 752 ± 0.05 | 763 ± 0.05 | 769 ± 0.05 | 765 ± 0.05 | 754 ± 0.05 | 760 ± 0.05 | 751 ± 0.05 | 761 ± 0.05 | 749 ± 0.05 | 750 ± 0.05 | 750 ± 0.05 | 756 ± 0.04 | 752 ± 0.05 | 757 ± 0.05 | 749 ± 0.05 | 767 ± 0.05 | 761 ± 0.06 | 758 ± 0.04 | |
| | .70 | 703 ± 0.05 | 716 ± 0.05 | 719 ± 0.06 | 716 ± 0.05 | 721 ± 0.05 | 717 ± 0.05 | 706 ± 0.05 | 722 ± 0.06 | 712 ± 0.05 | 711 ± 0.05 | 712 ± 0.05 | 712 ± 0.05 | 712 ± 0.05 | 715 ± 0.05 | 700 ± 0.05 | 692 ± 0.05 | 719 ± 0.06 | 709 ± 0.06 | 712 ± 0.05 |
| $MinCoeff$ | .99 | 1.003 ± 0.11 | 1.000 ± 0.11 | 1.001 ± 0.11 | 1.000 ± 0.11 | 1.000 ± 0.11 | 998 ± 0.11 | 1.002 ± 0.11 | 1.002 ± 0.11 | 998 ± 0.11 | 1.001 ± 0.11 | 1.003 ± 0.11 | 1.003 ± 0.11 | 1.001 ± 0.11 | 1.001 ± 0.11 | 998 ± 0.11 | 1.000 ± 0.11 | 1.004 ± 0.11 | 999 ± 0.11 | 1.002 ± 0.11 |
| | .95 | 1.021 ± 0.11 | 1.003 ± 0.11 | 1.004 ± 0.11 | 992 ± 0.11 | 1.002 ± 0.11 | 1.000 ± 0.11 | 1.010 ± 0.11 | 1.005 ± 0.11 | 997 ± 0.11 | 1.004 ± 0.11 | 1.008 ± 0.11 | 1.002 ± 0.11 | 999 ± 0.11 | 1.001 ± 0.11 | 1.006 ± 0.11 | 997 ± 0.11 | 1.013 ± 0.11 | 1.012 ± 0.11 | |
| | .90 | 1.037 ± 0.11 | 1.006 ± 0.12 | 1.007 ± 0.12 | 991 ± 0.11 | 1.009 ± 0.12 | 1.008 ± 0.11 | 1.009 ± 0.12 | 1.014 ± 0.11 | 993 ± 0.11 | 1.009 ± 0.12 | 1.017 ± 0.11 | 1.005 ± 0.12 | 1.003 ± 0.12 | 1.000 ± 0.11 | 1.007 ± 0.12 | 1.000 ± 0.11 | 1.024 ± 0.11 | 1.031 ± 0.12 | |
| | .85 | 1.049 ± 0.12 | 1.013 ± 0.12 | 1.013 ± 0.12 | 999 ± 0.12 | 1.009 ± 0.12 | 1.008 ± 0.12 | 1.018 ± 0.12 | 1.019 ± 0.12 | 1.005 ± 0.12 | 1.009 ± 0.12 | 1.027 ± 0.11 | 1.003 ± 0.12 | 1.006 ± 0.12 | 1.012 ± 0.12 | 1.004 ± 0.12 | 1.004 ± 0.11 | 1.031 ± 0.11 | 1.042 ± 0.12 | |
| | .80 | 1.060 ± 0.12 | 1.033 ± 0.13 | 1.037 ± 0.12 | 997 ± 0.12 | 1.020 ± 0.12 | 1.021 ± 0.12 | 1.021 ± 0.12 | 1.024 ± 0.12 | 1.006 ± 0.11 | 1.030 ± 0.12 | 1.035 ± 0.11 | 1.008 ± 0.13 | 1.007 ± 0.13 | 1.023 ± 0.13 | 1.007 ± 0.13 | 1.006 ± 0.12 | 1.045 ± 0.12 | 1.057 ± 0.12 | |
| | .75 | 1.066 ± 0.13 | 1.025 ± 0.13 | 1.024 ± 0.12 | 1.002 ± 0.12 | 1.036 ± 0.12 | 1.036 ± 0.12 | 1.029 ± 0.13 | 1.030 ± 0.12 | 1.013 ± 0.12 | 1.035 ± 0.13 | 1.036 ± 0.12 | 1.011 ± 0.13 | 1.010 ± 0.13 | 1.039 ± 0.14 | 1.005 ± 0.13 | 1.003 ± 0.13 | 1.062 ± 0.12 | 1.074 ± 0.12 | |
| | .70 | 1.067 ± 0.14 | 1.023 ± 0.13 | 1.028 ± 0.13 | 1.002 ± 0.14 | 1.050 ± 0.13 | 1.029 ± 0.13 | 1.034 ± 0.13 | 1.033 ± 0.12 | 1.008 ± 0.13 | 1.052 ± 0.13 | 1.042 ± 0.13 | 1.014 ± 0.13 | 1.013 ± 0.13 | 1.073 ± 0.14 | 1.007 ± 0.14 | 1.002 ± 0.13 | 1.073 ± 0.13 | 1.097 ± 0.12 | |

Table B13: Results for eye: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SeNet | SeNet+EM | SR | SAT+SR | SAT+EM+SR | SeNet+SR | SeNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugInAUC | |
|-----------------|-----|------------|------------|------------|------------|------------|------------|------------|------------|-------------------|-------------|------------|------------|------------|-------------------|-------------------|------------|------------|------------|------------|
| \widehat{Err} | .99 | 433 ± 0.13 | 441 ± 0.14 | 415 ± 0.12 | 419 ± 0.12 | 413 ± 0.13 | 425 ± 0.13 | 444 ± 0.14 | 414 ± 0.12 | 413 ± 0.13 | 414 ± 0.13 | 427 ± 0.13 | 421 ± 0.13 | 414 ± 0.12 | 411 ± 0.13 | 412 ± 0.13 | 426 ± 0.13 | 425 ± 0.13 | 427 ± 0.13 | |
| | .95 | 435 ± 0.13 | 436 ± 0.15 | 417 ± 0.12 | 414 ± 0.12 | 418 ± 0.12 | 420 ± 0.14 | 436 ± 0.15 | 412 ± 0.12 | 415 ± 0.12 | 412 ± 0.12 | 411 ± 0.13 | 439 ± 0.14 | 419 ± 0.13 | 414 ± 0.12 | 405 ± 0.13 | 413 ± 0.13 | 421 ± 0.13 | 421 ± 0.13 | 421 ± 0.13 |
| | .90 | 435 ± 0.14 | 433 ± 0.16 | 415 ± 0.12 | 422 ± 0.12 | 413 ± 0.13 | 413 ± 0.14 | 439 ± 0.15 | 411 ± 0.13 | 412 ± 0.12 | 406 ± 0.13 | 451 ± 0.14 | 412 ± 0.13 | 408 ± 0.12 | 398 ± 0.13 | 412 ± 0.13 | 417 ± 0.13 | 418 ± 0.14 | 412 ± 0.13 | |
| | .85 | 432 ± 0.14 | 429 ± 0.16 | 410 ± 0.12 | 413 ± 0.13 | 416 ± 0.14 | 407 ± 0.14 | 431 ± 0.15 | 407 ± 0.13 | 393 ± 0.12 | 406 ± 0.13 | 457 ± 0.15 | 407 ± 0.13 | 405 ± 0.12 | 400 ± 0.13 | 415 ± 0.13 | 410 ± 0.14 | 411 ± 0.15 | 409 ± 0.14 | |

DEEP NEURAL NETWORK BENCHMARKS FOR SELECTIVE CLASSIFICATION

Table B19: Results for house: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | ϵ | DG | SAT | SAT+EM | SeINet | SeINet+EM | SR | SAT+SR | SAT+EM+SR | SeINet+SR | SeINet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugInAUC | | |
|-----------------|------------|------------|-------------|------------|-------------|-------------|------------|-------------|------------|-------------|--------------|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|------------|
| \widehat{Err} | 99 | 217 ± .007 | 126 ± .007 | 130 ± .007 | 131 ± .006 | 133 ± .006 | 138 ± .006 | 124 ± .007 | 129 ± .007 | 128 ± .008 | 132 ± .006 | 124 ± .006 | 123 ± .006 | 128 ± .007 | 140 ± .006 | 143 ± .006 | 120 ± .006 | 124 ± .006 | 127 ± .006 | | |
| | 95 | 208 ± .007 | 114 ± .006 | 120 ± .006 | 122 ± .006 | 112 ± .006 | 117 ± .006 | 113 ± .006 | 116 ± .006 | 114 ± .006 | 117 ± .006 | 116 ± .006 | 116 ± .006 | 116 ± .006 | 131 ± .006 | 146 ± .006 | 105 ± .006 | 117 ± .006 | 120 ± .006 | | |
| | 90 | 190 ± .007 | 098 ± .006 | 106 ± .006 | 111 ± .006 | 112 ± .006 | 101 ± .006 | 109 ± .006 | 103 ± .005 | 107 ± .006 | 109 ± .006 | 109 ± .006 | 109 ± .006 | 109 ± .006 | 109 ± .006 | 123 ± .006 | 147 ± .007 | 094 ± .006 | 110 ± .006 | 112 ± .006 | |
| | 85 | 175 ± .007 | 081 ± .005 | 089 ± .005 | 093 ± .005 | 096 ± .006 | 083 ± .006 | 078 ± .005 | 085 ± .005 | 096 ± .005 | 096 ± .005 | 096 ± .005 | 096 ± .005 | 096 ± .005 | 113 ± .006 | 147 ± .007 | 079 ± .006 | 107 ± .006 | 101 ± .006 | 103 ± .005 | |
| | 80 | 158 ± .008 | 068 ± .005 | 078 ± .006 | 077 ± .006 | 080 ± .006 | 064 ± .005 | 067 ± .006 | 073 ± .006 | 077 ± .006 | 077 ± .006 | 077 ± .006 | 077 ± .006 | 077 ± .006 | 084 ± .005 | 102 ± .006 | 149 ± .008 | 065 ± .006 | 088 ± .005 | 084 ± .006 | |
| | 75 | 141 ± .008 | 057 ± .005 | 069 ± .006 | 062 ± .006 | 066 ± .006 | 057 ± .005 | 057 ± .005 | 071 ± .006 | 062 ± .006 | 064 ± .005 | 065 ± .005 | 065 ± .005 | 065 ± .005 | 065 ± .005 | 065 ± .006 | 148 ± .008 | 056 ± .006 | 075 ± .006 | 078 ± .006 | |
| | 70 | 141 ± .008 | 050 ± .005 | 060 ± .006 | 055 ± .006 | 052 ± .005 | 049 ± .005 | 050 ± .005 | 061 ± .006 | 052 ± .005 | 054 ± .005 | 055 ± .005 | 055 ± .005 | 055 ± .005 | 055 ± .006 | 149 ± .008 | 053 ± .005 | 063 ± .006 | 067 ± .006 | 067 ± .006 | |
| | 65 | 141 ± .008 | 048 ± .002 | 064 ± .006 | 060 ± .006 | 062 ± .002 | 063 ± .003 | 067 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 | 069 ± .002 |
| | 60 | 141 ± .008 | 046 ± .004 | 057 ± .004 | 052 ± .004 | 050 ± .005 | 047 ± .004 | 051 ± .004 | 053 ± .004 | 044 ± .004 | 048 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | 049 ± .004 | |
| | 55 | 141 ± .008 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | 040 ± .006 | |
| $\hat{\phi}$ | 99 | 892 ± .006 | 811 ± .007 | 849 ± .007 | 853 ± .007 | 844 ± .007 | 842 ± .006 | 837 ± .008 | 849 ± .006 | 844 ± .006 | 850 ± .007 | 848 ± .007 | 847 ± .007 | 842 ± .006 | 847 ± .007 | 841 ± .007 | 844 ± .006 | 806 ± .007 | 844 ± .006 | 806 ± .007 | |
| | 95 | 811 ± .007 | 843 ± .007 | 849 ± .007 | 853 ± .007 | 844 ± .007 | 842 ± .006 | 837 ± .008 | 849 ± .006 | 844 ± .006 | 850 ± .007 | 848 ± .007 | 847 ± .007 | 842 ± .006 | 847 ± .007 | 841 ± .007 | 844 ± .006 | 806 ± .007 | 844 ± .006 | 806 ± .007 | |
| | 90 | 779 ± .008 | 800 ± .006 | 788 ± .007 | 794 ± .008 | 795 ± .008 | 772 ± .008 | 786 ± .007 | 785 ± .008 | 800 ± .007 | 791 ± .008 | 807 ± .007 | 794 ± .008 | 786 ± .007 | 774 ± .007 | 787 ± .007 | 784 ± .008 | 789 ± .008 | 799 ± .008 | 802 ± .008 | |
| | 85 | 725 ± .008 | 740 ± .008 | 739 ± .009 | 742 ± .008 | 727 ± .009 | 728 ± .008 | 742 ± .009 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | 742 ± .008 | |
| | 80 | 678 ± .008 | 686 ± .009 | 688 ± .009 | 693 ± .008 | 677 ± .009 | 686 ± .009 | 685 ± .010 | 688 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | 683 ± .009 | |
| | 75 | 594 ± .019 | 1000 ± .019 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 70 | 594 ± .019 | 1000 ± .019 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 65 | 594 ± .019 | 1000 ± .019 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 60 | 594 ± .019 | 1000 ± .019 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 55 | 594 ± .019 | 1000 ± .019 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| $MinCoeff$ | 99 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 95 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 90 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 85 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 80 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 75 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 70 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 65 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 60 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |
| | 55 | 996 ± .018 | 999 ± .018 | 999 ± .018 | 1001 ± .019 | 1001 ± .018 | 993 ± .018 | 1002 ± .018 | 998 ± .018 | 1002 ± .018 | 1001 ± .019 | 998 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1001 ± .018 | 1001 ± .019 | 1002 ± .019 | |

Table B20: Results for indian: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | ϵ | DG | SAT | SAT+EM | SeINet | SeINet+EM | SR | SAT+SR | SAT+EM+SR | SeINet+SR | SeINet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross |
|-----------------|------------|-------------|------------------|------------------|-------------|------------------|------------------|------------------|-------------|------------------|--------------|------------------|------------------|-------------|-------------|------------------|-------------|
| \widehat{Err} | 99 | 050 ± .005 | .038 ± .004 | .041 ± .005 | .046 ± .005 | .131 ± .008 | 036 ± 004 | .038 ± .004 | .039 ± .005 | .041 ± .005 | .146 ± .008 | .040 ± .005 | .040 ± .005 | .055 ± .005 | .055 ± .006 | .084 ± .006 | .037 ± .005 |
| | 95 | .038 ± .005 | .029 ± .004 | .023 ± .004 | .031 ± .005 | .107 ± .008 | .023 ± .004 | .027 ± .004 | .023 ± .004 | .029 ± .004 | .123 ± .008 | .024 ± .004 | 017 ± 003 | .022 ± .005 | .022 ± .006 | .079 ± .006 | .027 ± .004 |
| | 90 | .032 ± .005 | .020 ± .004 | .012 ± .003 | .022 ± .004 | .099 ± .008 | .011 ± .005 | .017 ± .003 | .012 ± .003 | .016 ± .003 | .115 ± .008 | .012 ± .003 | 010 ± 003 | .034 ± .005 | .036 ± .006 | .079 ± .006 | .019 ± .004 |
| | 85 | .028 ± .005 | .012 ± .003 | .003 ± .002 | .014 ± .003 | .113 ± .008 | .007 ± .002 | .012 ± .003 | .008 ± .002 | .013 ± .003 | .103 ± .007 | 006 ± 002 | 006 ± 002 | .029 ± .005 | .056 ± .006 | .074 ± .006 | .013 ± .003 |
| | 80 | .024 ± .004 | .012 ± .003 | .005 ± .002 | .012 ± .003 | .094 ± .007 | .007 ± .002 | .012 ± .003 | .005 ± .002 | .031 ± .004 | .139 ± .009 | 003 ± 001 | .005 ± .002 | .022 ± .004 | .051 ± .006 | .072 ± .007 | .007 ± .002 |
| | 75 | .020 ± .004 | .011 ± .003 | .004 ± .001 | .007 ± .002 | .087 ± .008 | 003 ± 001 | .010 ± .003 | .004 ± .001 | .006 ± .002 | .125 ± .011 | 003 ± 001 | 003 ± 001 | .018 ± .004 | .044 ± .006 | .066 ± .006 | .006 ± .002 |
| | 70 | .016 ± .004 | .011 ± .003 | 002 ± 001 | .005 ± .002 | .068 ± .007 | 003 ± 002 | .010 ± .003 | .003 ± .001 | .005 ± .002 | .137 ± .011 | .003 ± .001 | .016 ± .004 | .043 ± .006 | .061 ± .006 | .004 ± .002 | |
| | 65 | .012 ± .003 | 000 ± 002 | 002 ± 002 | .002 ± .002 | .080 ± .005 | .082 ± .005 | 001 ± 002 | .086 ± .003 | .092 ± .002 | .150 ± .012 | .092 ± .002 | .095 ± .003 | .092 ± .002 | .092 ± .002 | .092 ± .002 | .092 ± .002 |
| | 60 | .010 ± .003 | .002 ± .002 | .002 ± .002 | .002 ± .002 | .080 ± .005 | .082 ± .005 | .088 ± .003 | .091 ± .002 | .096 ± .003 | .150 ± .012 | .092 ± .002 | .095 ± .003 | .092 ± .002 | .092 ± .002 | .092 ± .002 | .092 ± .002 |
| | 55 | .008 ± .003 | .002 ± .002 | .002 ± .002 | .002 ± .002 | .080 ± .005 | .082 ± .005 | .088 ± .003 | .091 ± .002 | .096 ± .003 | .150 ± .012 | .092 ± .002 | .095 ± .003 | .092 ± .002 | .092 ± .002 | .092 ± .002 | .092 ± .002 |
| $\hat{\phi}$ | 99 | 877 ± .008 | .007 ± .002 | .884 ± .008 | .896 ± .007 | 899 ± 006 | .885 ± .008 | .900 ± .007 | .891 ± .008 | .890 ± .008 | .888 ± .008 | .889 ± .008 | .882 ± .008 | .892 ± .007 | .908 ± .007 | 901 ± 019 | .919 ± .006 |
| | 95 | .828 ± .010 | .842 ± .008 | .848 ± .008 | .832 ± .008 | .858 ± .008 | .831 ± .010 | .836 ± .009 | .838 ± .009 | 849 ± 008 | .826 ± .010 | .838 ± .010 | .843 ± .009 | .839 ± .009 | .825 ± .009 | .869 ± .008 | |
| | 90 | .782 ± .010 | .796 ± .010 | .792 ± .010 | .805 ± .009 | 802 ± 009 | .784 ± .011 | .7 | | | | | | | | | |

Table B23: Results for letter: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | | | |
|-----------------|--------------|-------------|--------------------|--------------------|-------------------|-------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|-------------------|-------------------|------------|------------|
| \widehat{Err} | .99 | 122 ± .005 | .016 ± .002 | .014 ± .002 | .017 ± .002 | .026 ± .002 | .015 ± .002 | .014 ± .002 | .012 ± .002 | .014 ± .002 | .014 ± .002 | .014 ± .002 | .012 ± .002 | .012 ± .002 | .033 ± .003 | .061 ± .004 | .048 ± .003 | .020 ± .002 | | |
| | .95 | 105 ± .005 | .006 ± .001 | .005 ± .001 | .013 ± .002 | .012 ± .002 | .014 ± .002 | .006 ± .001 | .005 ± .001 | .006 ± .001 | .006 ± .001 | .006 ± .001 | .005 ± .001 | .005 ± .001 | .028 ± .003 | .058 ± .004 | .047 ± .003 | .009 ± .002 | | |
| | .90 | 084 ± .005 | .002 ± .001 | .002 ± .001 | .007 ± .001 | .009 ± .002 | .006 ± .001 | .002 ± .001 | .002 ± .001 | .002 ± .001 | .002 ± .001 | .002 ± .001 | .002 ± .001 | .002 ± .001 | .022 ± .003 | .048 ± .004 | .047 ± .004 | .005 ± .001 | | |
| | .85 | 070 ± .005 | .001 ± .001 | .002 ± .001 | .007 ± .001 | .007 ± .001 | .003 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .019 ± .002 | .041 ± .003 | .046 ± .004 | .002 ± .001 | | |
| | .80 | 053 ± .004 | .000 ± .000 | .002 ± .001 | .004 ± .001 | .006 ± .001 | .002 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .001 ± .001 | .001 ± .001 | .000 ± .000 | .015 ± .002 | .039 ± .003 | .044 ± .004 | .001 ± .001 | | |
| | .75 | .042 ± .004 | .000 ± .000 | .001 ± .001 | .003 ± .001 | .003 ± .001 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .014 ± .002 | .039 ± .003 | .044 ± .004 | .001 ± .001 | | |
| | .70 | .037 ± .004 | .000 ± .000 | .001 ± .001 | .003 ± .001 | .003 ± .001 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .001 ± .001 | .000 ± .000 | .012 ± .002 | .038 ± .003 | .042 ± .004 | .001 ± .001 | | |
| | $\hat{\phi}$ | .99 | 991 ± .001 | 994 ± .001 | 995 ± .001 | 990 ± .002 | 993 ± .001 | 994 ± .001 | 992 ± .002 | 993 ± .001 | 994 ± .001 | 992 ± .002 | 993 ± .001 | 996 ± .002 | 990 ± .002 | 993 ± .001 | 993 ± .001 | 990 ± .002 | 992 ± .002 | 991 ± .001 |
| | | .95 | 955 ± .003 | 957 ± .003 | 953 ± .004 | 956 ± .003 | 950 ± .003 | 957 ± .003 | 954 ± .003 | 954 ± .003 | 954 ± .003 | 954 ± .003 | 949 ± .004 | 949 ± .003 | 951 ± .003 | 952 ± .003 | 953 ± .003 | 954 ± .003 | 939 ± .004 | 959 ± .003 |
| | | .90 | 904 ± .004 | 903 ± .005 | 903 ± .005 | .911 ± .004 | 908 ± .004 | 907 ± .005 | 903 ± .005 | 905 ± .005 | 901 ± .005 | 901 ± .005 | 901 ± .005 | 902 ± .005 | 902 ± .005 | 903 ± .005 | 901 ± .005 | 907 ± .005 | 890 ± .005 | 923 ± .005 |
| .85 | | 860 ± .006 | 857 ± .006 | 848 ± .006 | 861 ± .005 | 853 ± .006 | 855 ± .006 | 850 ± .006 | 854 ± .006 | 851 ± .007 | 855 ± .006 | 854 ± .006 | 854 ± .006 | 854 ± .006 | 863 ± .005 | 849 ± .005 | 845 ± .005 | 882 ± .005 | 882 ± .005 | |
| .80 | | 807 ± .006 | .811 ± .006 | .801 ± .006 | .806 ± .006 | .810 ± .006 | .799 ± .007 | .810 ± .006 | .805 ± .006 | .798 ± .007 | .797 ± .007 | .797 ± .007 | .797 ± .007 | .811 ± .006 | .808 ± .006 | .803 ± .006 | .807 ± .006 | .789 ± .006 | 838 ± .006 | |
| .75 | | 756 ± .007 | .757 ± .007 | .751 ± .007 | .745 ± .007 | .751 ± .007 | .757 ± .007 | .764 ± .007 | .759 ± .007 | .749 ± .007 | .759 ± .007 | .753 ± .007 | .763 ± .007 | .759 ± .007 | .751 ± .007 | .766 ± .007 | .736 ± .006 | 793 ± .006 | 793 ± .006 | |
| .70 | | 707 ± .008 | .698 ± .007 | .703 ± .008 | .701 ± .007 | .711 ± .008 | .698 ± .008 | .702 ± .007 | .697 ± .007 | .708 ± .008 | .712 ± .008 | .716 ± .008 | .715 ± .007 | .701 ± .007 | .710 ± .007 | .690 ± .006 | .747 ± .007 | | | |

Table B24: Results for magic: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugNAUC | |
|-----------------|--------------|---------------|-------------------|--------------------|----------------------|-------------------|-------------------|--------------------|---------------|-------------------|---------------|----------------------|---------------|---------------|-------------------|---------------|---------------|---------------|---------------|---------------|
| \widehat{Err} | .99 | 225 ± .008 | 142 ± .007 | 155 ± .007 | 143 ± .007 | 148 ± .007 | 155 ± .008 | 143 ± .007 | 153 ± .007 | 141 ± .007 | 147 ± .008 | 153 ± .008 | 149 ± .007 | 147 ± .008 | 152 ± .008 | 157 ± .008 | 158 ± .008 | 159 ± .008 | 155 ± .008 | |
| | .95 | 217 ± .008 | 127 ± .007 | 141 ± .007 | 129 ± .007 | 144 ± .007 | 136 ± .007 | 128 ± .007 | 142 ± .007 | 127 ± .007 | 142 ± .007 | 147 ± .008 | 134 ± .007 | 134 ± .007 | 147 ± .008 | 153 ± .007 | 141 ± .007 | 153 ± .007 | 153 ± .008 | 146 ± .008 |
| | .90 | 203 ± .008 | 109 ± .007 | 123 ± .007 | 112 ± .007 | 125 ± .007 | 118 ± .007 | 110 ± .007 | 123 ± .007 | 108 ± .007 | 125 ± .007 | 140 ± .007 | 118 ± .007 | 120 ± .007 | 134 ± .008 | 139 ± .007 | 114 ± .007 | 125 ± .007 | 147 ± .008 | 146 ± .008 |
| | .85 | 186 ± .008 | 099 ± .007 | 108 ± .006 | 099 ± .007 | 121 ± .007 | 112 ± .007 | 097 ± .007 | 109 ± .006 | 098 ± .007 | 121 ± .007 | 134 ± .007 | 105 ± .007 | 107 ± .006 | 120 ± .008 | 130 ± .007 | 101 ± .007 | 124 ± .008 | 133 ± .008 | 133 ± .008 |
| | .80 | 170 ± .008 | 087 ± .006 | 092 ± .006 | 085 ± .006 | 100 ± .006 | 088 ± .006 | 088 ± .006 | 094 ± .006 | 086 ± .006 | 099 ± .006 | 120 ± .008 | 105 ± .006 | 099 ± .006 | 120 ± .008 | 130 ± .007 | 101 ± .007 | 124 ± .008 | 133 ± .008 | 133 ± .008 |
| | .75 | 156 ± .007 | 076 ± .006 | 084 ± .006 | 074 ± .006 | 080 ± .006 | 088 ± .006 | 073 ± .006 | 083 ± .006 | 075 ± .006 | 090 ± .006 | 119 ± .007 | 096 ± .006 | 088 ± .006 | 112 ± .007 | 124 ± .007 | 089 ± .006 | 108 ± .006 | 116 ± .006 | 116 ± .006 |
| | .70 | 144 ± .008 | 072 ± .006 | 072 ± .006 | 069 ± .006 | 086 ± .006 | 085 ± .006 | 073 ± .006 | 070 ± .006 | 073 ± .006 | 087 ± .006 | 105 ± .007 | 072 ± .006 | 072 ± .006 | 105 ± .007 | 124 ± .007 | 078 ± .006 | 093 ± .007 | 103 ± .008 | 103 ± .008 |
| | $\hat{\phi}$ | .99 | 986 ± .003 | 981 ± .002 | 983 ± .001 | 990 ± .002 | 988 ± .002 | 997 ± .001 | 987 ± .002 | 992 ± .002 | 989 ± .002 | 990 ± .002 | 984 ± .002 | 985 ± .002 | 990 ± .002 | 992 ± .002 | 992 ± .002 | 992 ± .002 | 991 ± .002 | 991 ± .002 |
| | | .95 | 954 ± .004 | 936 ± .005 | 953 ± .005 | 943 ± .004 | 959 ± .004 | 945 ± .004 | 942 ± .005 | 959 ± .004 | 951 ± .005 | 952 ± .004 | 953 ± .005 | 941 ± .005 | 936 ± .005 | 957 ± .004 | 956 ± .004 | 951 ± .004 | 958 ± .004 | 949 ± .005 |
| | | .90 | 904 ± .006 | 880 ± .006 | 891 ± .007 | 896 ± .006 | 897 ± .006 | 883 ± .006 | 890 ± .006 | 892 ± .006 | 887 ± .007 | 898 ± .006 | 893 ± .006 | 885 ± .006 | 889 ± .006 | 909 ± .006 | 887 ± .006 | 887 ± .006 | 904 ± .006 | 907 ± .006 |
| .85 | | 849 ± .006 | 843 ± .007 | 841 ± .008 | 847 ± .007 | 855 ± .007 | 850 ± .007 | 847 ± .007 | 849 ± .007 | 851 ± .007 | 870 ± .007 | 857 ± .007 | 845 ± .007 | 847 ± .007 | 857 ± .007 | 845 ± .007 | 860 ± .007 | 864 ± .007 | 870 ± .007 | |
| .80 | | 798 ± .007 | 784 ± .008 | 791 ± .008 | 788 ± .008 | 788 ± .008 | 788 ± .008 | 778 ± .008 | 787 ± .008 | 790 ± .008 | 809 ± .008 | 819 ± .009 | 793 ± .008 | 799 ± .007 | 800 ± .007 | 805 ± .007 | 805 ± .007 | 817 ± .008 | 822 ± .007 | |
| .75 | | 748 ± .008 | 724 ± .009 | 744 ± .008 | 734 ± .008 | 762 ± .009 | 736 ± .009 | 734 ± .009 | 744 ± .008 | 734 ± .008 | 760 ± .009 | 776 ± .009 | 745 ± .009 | 747 ± .008 | 752 ± .008 | 760 ± .008 | 762 ± .008 | 784 ± .008 | 798 ± .008 | |
| .70 | | 698 ± .008 | 687 ± .010 | 698 ± .009 | 688 ± .009 | 700 ± .009 | 693 ± .010 | 686 ± .009 | 697 ± .008 | 706 ± .008 | 705 ± .009 | 700 ± .009 | 694 ± .009 | 704 ± .009 | 701 ± .009 | 719 ± .009 | 719 ± .009 | 714 ± .010 | 719 ± .009 | |
| $MinCoeff$ | | .99 | 1.001 ± 0.019 | 996 ± 0.019 | 997 ± 0.019 | 1.001 ± 0.019 | 996 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 1.002 ± 0.019 | 1.003 ± 0.019 |
| | | .95 | 1.007 ± 0.020 | 989 ± 0.019 | 1.000 ± 0.019 | 1.002 ± 0.019 | 996 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 998 ± 0.019 | 998 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 999 ± 0.019 | 1.002 ± 0.020 | 1.029 ± 0.020 |
| | | .90 | 1.017 ± 0.020 | 994 ± 0.020 | 998 ± 0.019 | 999 ± 0.019 | 998 ± 0.020 | 1.000 ± 0.020 | 999 ± 0.020 | 998 ± 0.019 | 994 ± 0.020 | 1.001 ± 0.019 | 1.003 ± 0.020 | 999 ± 0.020 | 1.034 ± 0.021 | 1.034 ± 0.021 | 989 ± 0.020 | 1.022 ± 0.021 | 1.046 ± 0.021 | 1.046 ± 0.021 |
| | .85 | 1.023 ± 0.021 | 997 ± 0.022 | 995 ± 0.021 | 1.010 ± 0.021 | 1.003 ± 0.021 | 1.027 ± 0.021 | 999 ± 0.021 | 995 ± 0.021 | 1.009 ± 0.020 | 1.005 ± 0.021 | 1.009 ± 0.020 | 1.005 ± 0.021 | 1.003 ± 0.021 | 1.007 ± 0.022 | 1.055 ± 0.021 | 1.024 ± 0.020 | 985 ± 0.021 | 1.060 ± 0.022 | |
| | .80 | 1.034 ± 0.022 | 1.010 ± 0.023 | 999 ± 0.021 | 1.017 ± 0.022 | 1.017 ± 0.022 | 1.017 ± 0.022 | 1.009 ± 0.023 | 1.005 ± 0.022 | 1.009 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | 1.016 ± 0.021 | |
| | .75 | 1.044 ± 0.022 | 1.023 ± 0.024 | 1.014 ± 0.024 | 1.034 ± 0.024 | 1.027 ± 0.023 | 1.057 ± 0.023 | 1.018 ± 0.023 | 1.024 ± 0.023 | 1.025 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | 1.024 ± 0.023 | |
| | .70 | 1.066 ± 0.022 | 1.044 ± 0.024 | 1.038 ± 0.025 | 1.057 ± 0.024 | 1.057 ± 0.024 | 1.065 ± 0.023 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | 1.048 ± 0.024 | |

Table B25: Results for miniboone: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PlugNAUC |
|-----------------|-----|------------|-------------|--------------------|-------------|-------------|-------------|--------------------|--------------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| \widehat{Err} | .99 | 065 ± .002 | .061 ± .002 | .057 ± .002 | .059 ± .002 | .066 ± .002 | .061 ± .002 | .060 ± .002 | .058 ± .002 | .060 ± .002 | .067 ± .002 | .062 ± .002 | .060 ± .002 | .067 ± .002 | .075 ± .002 | .073 ± .002 | .073 ± .002 | .068 ± .002 | .062 ± .002 |
| | .95 | 055 ± .002 | .049 ± .002 | .047 ± .002 | .047 ± .002 | .052 ± .002 | .046 ± .002 | .044 ± .002 | .044 ± .002 | .046 ± .002 | .050 ± .002 | .051 ± .002 | .046 ± .002 | .046 ± .002 | .051 ± .002 | .051 ± .002 | .049 ± .002 | .046 ± .002 | .046 ± .002 |
| | .90 | 047 ± .002 | .036 ± .002 | .033 ± .002 | .033 ± .002 | .038 ± .002 | .033 ± .002 | .032 ± .002 | .032 ± .002 | .033 ± .002 | .035 ± .002 | .035 ± .002 | .030 ± .002 | .034 ± .002 | .038 ± .002 | .037 ± .002 | .033 ± .002 | .036 ± .002 | .036 ± .002 |
| | .85 | 041 ± .002 | .024 ± .001 | .023 ± .001 | .024 ± .001 | .027 ± .001 | .023 ± .001 | .023 ± .001 | .023 ± .001 | .024 ± .001 | .029 ± .001 | .031 ± .001 | .023 ± .001 | .027 ± .001 | .035 ± .002 | .035 ± .002 | .035 ± .002 | .035 ± .002 | .035 ± .002 |
| | .80 | 036 ± .002 | .017 ± .001 | .017 ± .001 | .018 ± .001 | .020 ± .001 | .018 ± .001 | .018 ± .001 | .018 ± .001 | .018 ± .001 | .020 ± .001 | .023 ± .001 | .0 | | | | | | |

DEEP NEURAL NETWORK BENCHMARKS FOR SELECTIVE CLASSIFICATION

Table B28: Results for online: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhgnAUC | | |
|-----------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|---------------|---------------|---------------|---------------|---------------|---------------|-------------|
| \widehat{Err} | 99 | 061 ± 0.006 | 055 ± 0.006 | 033 ± 0.006 | 031 ± 0.006 | 099 ± 0.006 | 084 ± 0.006 | 084 ± 0.006 | 082 ± 0.006 | 102 ± 0.006 | 101 ± 0.006 | 089 ± 0.006 | 089 ± 0.006 | 082 ± 0.006 | 084 ± 0.006 | 082 ± 0.006 | 083 ± 0.006 | 086 ± 0.006 | 089 ± 0.006 | | |
| | 90 | 076 ± 0.006 | 087 ± 0.006 | 078 ± 0.006 | 083 ± 0.006 | 089 ± 0.006 | 081 ± 0.006 | 081 ± 0.006 | 082 ± 0.006 | 084 ± 0.006 | 088 ± 0.006 | 091 ± 0.006 | 078 ± 0.006 | 081 ± 0.006 | 083 ± 0.006 | 106 ± 0.007 | 085 ± 0.006 | 095 ± 0.006 | 099 ± 0.006 | | |
| | 90 | 066 ± 0.005 | 073 ± 0.005 | 067 ± 0.005 | 071 ± 0.005 | 068 ± 0.006 | 064 ± 0.005 | 067 ± 0.005 | 065 ± 0.005 | 068 ± 0.005 | 069 ± 0.006 | 084 ± 0.006 | 063 ± 0.005 | 069 ± 0.005 | 087 ± 0.006 | 108 ± 0.007 | 070 ± 0.005 | 095 ± 0.006 | 098 ± 0.006 | | |
| | 85 | 056 ± 0.005 | 061 ± 0.005 | 063 ± 0.005 | 069 ± 0.006 | 062 ± 0.006 | 051 ± 0.005 | 053 ± 0.004 | 055 ± 0.005 | 060 ± 0.006 | 062 ± 0.006 | 075 ± 0.006 | 049 ± 0.005 | 051 ± 0.005 | 082 ± 0.006 | 111 ± 0.008 | 047 ± 0.005 | 088 ± 0.006 | 095 ± 0.007 | | |
| | 80 | 048 ± 0.005 | 047 ± 0.005 | 040 ± 0.004 | 046 ± 0.005 | 051 ± 0.005 | 040 ± 0.005 | 042 ± 0.004 | 042 ± 0.004 | 047 ± 0.005 | 049 ± 0.005 | 068 ± 0.005 | 038 ± 0.004 | 043 ± 0.005 | 079 ± 0.006 | 113 ± 0.008 | 039 ± 0.004 | 081 ± 0.006 | 094 ± 0.007 | | |
| | 75 | 039 ± 0.005 | 038 ± 0.004 | 031 ± 0.004 | 041 ± 0.005 | 041 ± 0.005 | 032 ± 0.004 | 035 ± 0.004 | 031 ± 0.004 | 044 ± 0.005 | 042 ± 0.004 | 060 ± 0.005 | 033 ± 0.004 | 036 ± 0.004 | 071 ± 0.006 | 115 ± 0.008 | 032 ± 0.004 | 087 ± 0.006 | 084 ± 0.007 | | |
| | 70 | 032 ± 0.004 | 031 ± 0.004 | 026 ± 0.004 | 030 ± 0.004 | 039 ± 0.005 | 028 ± 0.004 | 027 ± 0.004 | 028 ± 0.004 | 033 ± 0.004 | 035 ± 0.004 | 052 ± 0.005 | 028 ± 0.004 | 033 ± 0.005 | 057 ± 0.005 | 120 ± 0.008 | 028 ± 0.004 | 065 ± 0.006 | 071 ± 0.007 | | |
| | $\hat{\phi}$ | 99 | 988 ± 0.002 | 988 ± 0.002 | 991 ± 0.002 | 991 ± 0.002 | 988 ± 0.002 | 985 ± 0.002 | 990 ± 0.002 | 987 ± 0.002 | 992 ± 0.002 | 990 ± 0.002 | 994 ± 0.002 | 987 ± 0.002 | 986 ± 0.002 | 989 ± 0.002 | 993 ± 0.002 | 993 ± 0.002 | 993 ± 0.002 | 998 ± 0.002 | |
| | | 95 | 948 ± 0.005 | 955 ± 0.006 | 953 ± 0.006 | 951 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 | 954 ± 0.006 |
| | | 90 | 902 ± 0.006 | 907 ± 0.006 | 925 ± 0.006 | 908 ± 0.006 | 920 ± 0.006 | 914 ± 0.006 | 915 ± 0.005 | 920 ± 0.006 | 901 ± 0.007 | 923 ± 0.006 | 896 ± 0.006 | 913 ± 0.006 | 916 ± 0.006 | 895 ± 0.006 | 934 ± 0.007 | 934 ± 0.007 | 916 ± 0.006 | 902 ± 0.005 | |
| 85 | | 809 ± 0.008 | 868 ± 0.007 | 809 ± 0.008 | 869 ± 0.007 | 868 ± 0.007 | 862 ± 0.008 | 868 ± 0.007 | 868 ± 0.007 | 868 ± 0.007 | 870 ± 0.007 | 871 ± 0.007 | 851 ± 0.007 | 864 ± 0.008 | 857 ± 0.008 | 834 ± 0.007 | 829 ± 0.008 | 854 ± 0.008 | 858 ± 0.007 | 844 ± 0.007 | |
| 80 | | 829 ± 0.008 | 821 ± 0.009 | 820 ± 0.009 | 820 ± 0.008 | 822 ± 0.009 | 819 ± 0.009 | 819 ± 0.008 | 822 ± 0.009 | 820 ± 0.008 | 824 ± 0.009 | 816 ± 0.008 | 820 ± 0.009 | 821 ± 0.008 | 790 ± 0.008 | 795 ± 0.009 | 806 ± 0.008 | 799 ± 0.008 | 790 ± 0.008 | | |
| 75 | | 779 ± 0.009 | 777 ± 0.009 | 767 ± 0.009 | 775 ± 0.010 | 776 ± 0.009 | 779 ± 0.009 | 776 ± 0.009 | 776 ± 0.009 | 769 ± 0.009 | 780 ± 0.009 | 779 ± 0.009 | 768 ± 0.009 | 772 ± 0.009 | 778 ± 0.008 | 783 ± 0.009 | 766 ± 0.010 | 740 ± 0.009 | 731 ± 0.008 | | |
| 70 | | 708 ± 0.009 | 732 ± 0.009 | 712 ± 0.009 | 721 ± 0.011 | 729 ± 0.010 | 732 ± 0.009 | 716 ± 0.010 | 719 ± 0.009 | 718 ± 0.010 | 718 ± 0.010 | 707 ± 0.011 | 729 ± 0.009 | 710 ± 0.009 | 680 ± 0.009 | 690 ± 0.009 | 620 ± 0.010 | 608 ± 0.010 | 675 ± 0.009 | | |
| $MinCoeff$ | | 99 | 983 ± 0.046 | 991 ± 0.047 | 990 ± 0.049 | 978 ± 0.047 | 976 ± 0.048 | 972 ± 0.047 | 982 ± 0.048 | 968 ± 0.048 | 972 ± 0.047 | 983 ± 0.047 | 994 ± 0.047 | 994 ± 0.048 | 981 ± 0.048 | 991 ± 0.048 | 1.015 ± 0.048 | 993 ± 0.048 | 1.006 ± 0.048 | 1.008 ± 0.047 | |
| | | 95 | 913 ± 0.047 | 933 ± 0.045 | 887 ± 0.046 | 886 ± 0.046 | 890 ± 0.048 | 911 ± 0.050 | 907 ± 0.047 | 901 ± 0.047 | 911 ± 0.046 | 900 ± 0.046 | 933 ± 0.046 | 901 ± 0.047 | 900 ± 0.048 | 986 ± 0.050 | 1.000 ± 0.049 | 938 ± 0.047 | 1.010 ± 0.046 | 1.018 ± 0.047 | |
| | | 90 | 811 ± 0.046 | 759 ± 0.042 | 822 ± 0.047 | 751 ± 0.047 | 782 ± 0.046 | 816 ± 0.048 | 818 ± 0.046 | 809 ± 0.046 | 846 ± 0.046 | 734 ± 0.044 | 803 ± 0.045 | 821 ± 0.048 | 780 ± 0.046 | 793 ± 0.044 | 972 ± 0.051 | 1.045 ± 0.051 | 867 ± 0.046 | 1.031 ± 0.048 | |
| | 85 | 735 ± 0.045 | 688 ± 0.039 | 675 ± 0.045 | 727 ± 0.042 | 699 ± 0.043 | 630 ± 0.042 | 686 ± 0.042 | 667 ± 0.044 | 721 ± 0.041 | 700 ± 0.043 | 750 ± 0.045 | 626 ± 0.042 | 657 ± 0.042 | 945 ± 0.052 | 1.046 ± 0.052 | 630 ± 0.044 | 1.012 ± 0.052 | 1.046 ± 0.052 | | |
| | 80 | 684 ± 0.045 | 592 ± 0.038 | 585 ± 0.045 | 472 ± 0.038 | 459 ± 0.038 | 454 ± 0.043 | 567 ± 0.041 | 590 ± 0.044 | 536 ± 0.038 | 489 ± 0.039 | 685 ± 0.043 | 533 ± 0.040 | 548 ± 0.041 | 900 ± 0.053 | 1.063 ± 0.051 | 539 ± 0.041 | 997 ± 0.055 | 1.042 ± 0.056 | | |
| | 75 | 614 ± 0.047 | 374 ± 0.033 | 484 ± 0.040 | 264 ± 0.031 | 267 ± 0.030 | 467 ± 0.040 | 492 ± 0.039 | 479 ± 0.040 | 394 ± 0.031 | 270 ± 0.031 | 605 ± 0.043 | 462 ± 0.039 | 411 ± 0.039 | 896 ± 0.050 | 1.064 ± 0.055 | 475 ± 0.039 | 965 ± 0.058 | 1.040 ± 0.057 | | |
| | 70 | 580 ± 0.046 | 205 ± 0.027 | 416 ± 0.039 | 193 ± 0.028 | 250 ± 0.030 | 418 ± 0.040 | 430 ± 0.038 | 426 ± 0.038 | 229 ± 0.029 | 489 ± 0.043 | 372 ± 0.039 | 363 ± 0.040 | 730 ± 0.053 | 1.094 ± 0.055 | 433 ± 0.037 | 964 ± 0.060 | 994 ± 0.058 | | | |

Table B29: Results for orgasmnist: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | | |
|-----------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| \widehat{Err} | 99 | 005 ± 0.001 | 002 ± 0.001 | 003 ± 0.001 | 003 ± 0.001 | 001 ± 0.000 | 002 ± 0.000 | 002 ± 0.000 | 002 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 002 ± 0.000 | | |
| | 95 | 002 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 001 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 002 ± 0.001 | 008 ± 0.001 | 009 ± 0.001 | |
| | 90 | 001 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 008 ± 0.001 | 009 ± 0.001 | |
| | 85 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 008 ± 0.001 | 009 ± 0.001 | |
| | 80 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 008 ± 0.001 | 009 ± 0.001 | |
| | 75 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 008 ± 0.001 | 009 ± 0.001 | |
| | 70 | 001 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 000 ± 0.000 | 001 ± 0.000 | 008 ± 0.001 | 009 ± 0.001 | |
| | $\hat{\phi}$ | 99 | 987 ± 0.001 | 990 ± 0.001 | 989 ± 0.001 | 987 ± 0.001 | 989 ± 0.001 | 991 ± 0.001 | 990 ± 0.001 | 987 ± 0.001 | 992 ± 0.001 | 990 ± 0.001 | 994 ± 0.001 | 990 ± 0.001 | 988 ± 0.001 | 991 ± 0.001 | 991 ± 0.001 | 990 ± 0.001 | 990 ± 0.001 |
| | | 95 | 952 ± 0.002 | 954 ± 0.002 | 950 ± 0.002 | 951 ± 0.002 | 955 ± 0.002 | 954 ± 0.002 | 950 ± 0.002 | 949 ± 0.002 | 954 ± 0.002 | 954 ± 0.002 | 956 ± 0.002 | 955 ± 0.002 | 952 ± 0.002 | 952 ± 0.002 | 952 ± 0.002 | 952 ± 0.002 | 956 ± 0.002 |
| | | 90 | 905 ± 0.002 | 904 ± 0.002 | 906 ± 0.002 | 907 ± 0.002 | 905 ± 0.002 | 904 ± 0.002 | 901 ± 0.002 | 901 ± 0.002 | 907 ± 0.002 | 903 ± 0.002 | 904 ± 0.002 | 903 ± 0.002 | 905 ± 0.002 | 905 ± 0.002 | 905 ± 0.002 | 905 ± 0.002 | 905 ± 0.002 |
| 85 | | 848 ± 0.003 | 853 ± 0.003 | 849 ± 0.003 | 853 ± 0.003 | 853 ± 0.003 | 859 ± 0.003 | 852 ± 0.003 | 846 ± 0.003 | 858 ± 0.003 | 858 ± 0.003 | 858 ± 0.003 | 853 ± 0.003 | 853 ± 0.003 | 855 ± 0.003 | 855 ± 0.003 | 855 ± 0.003 | 859 ± 0.003 | |
| 80 | | 800 ± 0.004 | 805 ± 0.003 | 802 ± 0.003 | 807 ± 0.004 | 806 ± 0.004 | 806 ± 0.004 | 808 ± 0.003 | 802 ± 0.004 | 808 ± 0.004 | 808 ± 0.004 | 805 ± 0.002 | 805 ± 0.003 | 805 ± 0.003 | 807 ± 0.003 | 803 ± 0.004 | 805 ± 0.004 | 829 ± 0.003 | |
| 75 | | 750 ± 0.003 | 755 ± 0.004 | 755 ± 0.004 | 768 ± 0.004 | 760 ± 0.004 | 758 ± 0.004 | 757 ± 0.004 | 758 ± 0.004 | 758 ± 0.004 | 758 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | 753 ± 0.004 | |
| 70 | | 705 ± 0.004 | 698 ± 0.004 | 700 ± 0.004 | 710 ± 0.004 | 709 ± 0.004 | 706 ± 0.004 | 720 ± 0.004 | 703 ± 0.004 | 715 ± 0.004 | 715 ± 0.004 | 703 ± 0.004 | 712 ± 0.004 | 710 ± 0.004 | 709 ± 0.004 | 705 ± 0.004 | 700 ± 0.004 | 745 ± 0.003 | |

Table B30: Results for orgasmnist: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross |
|-----------------|----|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|
| \widehat{Err} | 99 | 034 ± 0.003 | 023 ± 0.002 | 023 ± 0.002 | 021 ± 0.002 | 025 ± 0.002 | 020 ± 0.002 | 023 ± 0.002 | 021 ± 0.002 | 020 ± 0.002 | 022 ± 0.002 | 019 ± 0.002 | 018 ± 0.002 | 033 ± 0.002 | 043 ± 0.003 | 046 ± 0.003 | 017 ± 0.002 |
| | 95 | 023 ± 0.002 | 012 ± 0.002 | 014 ± 0.002 | 013 ± 0.002 | 016 ± 0.002 | 010 ± 0.001 | 010 ± 0.001 | 012 ± 0.002 | 009 ± 0.001 | 012 ± 0.002 | 008 ± 0.001 | 005 ± 0.001 | 027 ± 0.002 | 039 ± 0.003 | 045 ± 0.003 | 007 ± 0.001 |
| | 90 | 015 ± 0.002 | 004 ± 0.001 | 007 ± 0.001 | 007 ± 0.001 | 005 ± 0.001 | 003 ± 0.001 | 003 ± 0.001 | 001 ± 0.001 | 001 ± 0.001 | 002 ± 0.001 | 003 ± 0.001 | 003 ± 0.001 | 005 ± 0.001 | 035 ± 0.003 | 047 ± 0.003 | 003 ± 0.001 |
| | 85 | 011 ± 0.002 | 002 ± 0.001 | 004 ± 0.001 | 006 ± 0.001 | 003 ± 0.001 | 002 ± 0.001 | | | | | | | | | | |

Table B33: Results for pathmnist: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfiNet | SELE | REG | SCross |
|-----------------|-----|-------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| \widehat{Err} | .99 | .027 ± .001 | .011 ± .001 | .020 ± .001 | .014 ± .001 | .022 ± .001 | .010 ± .001 | .011 ± .001 | .019 ± .001 | .012 ± .001 | .019 ± .001 | .007 ± .001 | 0.06 ± 0.01 | .020 ± .001 | .028 ± .001 | .033 ± .001 | .008 ± .001 |
| | .95 | .019 ± .001 | .004 ± .000 | .010 ± .001 | .007 ± .001 | .016 ± .001 | .003 ± .000 | .003 ± .000 | .008 ± .001 | .005 ± .001 | .013 ± .001 | .002 ± .000 | 0.01 ± 0.00 | .006 ± .001 | .028 ± .001 | .031 ± .001 | .000 ± .000 |
| | .90 | .013 ± .001 | .002 ± .000 | .004 ± .001 | .003 ± .000 | .023 ± .001 | .001 ± .000 | .002 ± .000 | .004 ± .000 | .002 ± .000 | .022 ± .001 | .001 ± .000 | 0.00 ± 0.00 | .002 ± .000 | .027 ± .001 | .035 ± .001 | .001 ± .000 |
| | .85 | .010 ± .001 | .001 ± .000 | .002 ± .000 | .001 ± .000 | .004 ± .000 | .001 ± .000 | .001 ± .000 | .002 ± .000 | .001 ± .000 | .025 ± .001 | .000 ± .000 | 0.00 ± 0.00 | .000 ± .000 | .025 ± .001 | .036 ± .001 | .001 ± .000 |
| | .80 | .007 ± .001 | 0.00 ± 0.00 | .001 ± .000 | .001 ± .000 | .002 ± .000 | .001 ± .000 | 0.00 ± 0.00 | .001 ± .000 | .001 ± .000 | .031 ± .000 | .000 ± .000 | 0.00 ± 0.00 | .000 ± .000 | .024 ± .001 | .037 ± .001 | .001 ± .000 |
| | .75 | .006 ± .001 | 0.00 ± 0.00 | .000 ± .000 | 0.00 ± 0.00 | .012 ± .001 | .001 ± .000 | 0.00 ± 0.00 | .001 ± .000 | 0.00 ± 0.00 | .010 ± .001 | 0.00 ± 0.00 | 0.00 ± 0.00 | .000 ± .000 | .023 ± .001 | .038 ± .001 | 0.00 ± 0.00 |
| | .70 | .005 ± .001 | 0.07 ± 0.00 | .001 ± .000 | .001 ± .000 | .001 ± .000 | .001 ± .000 | 0.00 ± 0.00 | .001 ± .000 | 0.00 ± 0.00 | .001 ± .000 | 0.00 ± 0.00 | 0.00 ± 0.00 | .001 ± .000 | .022 ± .001 | .038 ± .001 | 0.00 ± 0.00 |
| $\hat{\phi}$ | .99 | .989 ± .001 | .989 ± .001 | 0.99 ± 0.01 | .990 ± .001 | .989 ± .001 | .991 ± .001 | 0.99 ± 0.01 | .989 ± .001 | .989 ± .001 | 0.99 ± 0.01 | .989 ± .001 | .989 ± .001 | .989 ± .001 | .992 ± .001 | 0.99 ± 0.01 | .992 ± .001 |
| | .95 | .948 ± .002 | 0.95 ± 0.02 | 0.94 ± 0.01 | .951 ± .002 | .951 ± .002 | 0.95 ± 0.02 | 0.95 ± 0.01 | 0.95 ± 0.01 | .951 ± .002 | .948 ± .002 | .951 ± .002 | .951 ± .002 | .947 ± .001 | .955 ± .002 | .944 ± .002 | .958 ± .001 |
| | .90 | .898 ± .002 | 0.91 ± 0.02 | 0.89 ± 0.02 | 0.89 ± 0.02 | .899 ± .002 | 0.90 ± 0.02 | .902 ± .002 | 0.91 ± 0.02 | .901 ± .002 | .898 ± .002 | .900 ± .002 | .902 ± .002 | .895 ± .002 | .902 ± .002 | .895 ± .002 | .914 ± .002 |
| | .85 | .849 ± .002 | .848 ± .002 | .852 ± .003 | .852 ± .002 | .853 ± .002 | 0.85 ± 0.02 | .848 ± .002 | .852 ± .002 | .854 ± .002 | .856 ± .002 | .851 ± .002 | .852 ± .002 | .845 ± .002 | .851 ± .002 | .847 ± .002 | .863 ± .002 |
| | .80 | .797 ± .003 | .798 ± .002 | .801 ± .003 | .807 ± .002 | .805 ± .003 | 0.80 ± 0.03 | .798 ± .003 | .804 ± .003 | .803 ± .002 | .799 ± .002 | .796 ± .002 | .793 ± .003 | .799 ± .003 | 0.80 ± 0.02 | .807 ± .003 | |
| | .75 | .749 ± .003 | 0.75 ± 0.03 | .754 ± .003 | .754 ± .003 | .745 ± .003 | .744 ± .003 | .752 ± .003 | .752 ± .003 | 0.75 ± 0.03 | .752 ± .003 | .751 ± .003 | .744 ± .003 | .742 ± .003 | 0.75 ± 0.03 | .757 ± .003 | .748 ± .003 |
| | .70 | .697 ± .003 | .703 ± .003 | .702 ± .003 | .706 ± .003 | 0.70 ± 0.03 | .693 ± .003 | .704 ± .003 | .706 ± .003 | .701 ± .003 | .705 ± .003 | .695 ± .003 | .690 ± .003 | 0.71 ± 0.03 | .701 ± .003 | .701 ± .003 | .695 ± .003 |

Table B34: Results for pathmnist: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfiNet | SELE | REG | SCross | AUCross | PlugInAUC |
|-----------------|-----|--------------|--------------|--------------|---------------------|---------------------|--------------------|--------------|--------------------|--------------------|--------------|---------------------|---------------------|--------------------|--------------------|--------------------|--------------|---------------------|--------------------|
| \widehat{Err} | .99 | .255 ± .017 | .133 ± .013 | .141 ± .015 | .134 ± .014 | .161 ± .015 | .137 ± .015 | .136 ± .014 | .142 ± .015 | 0.17 ± 0.03 | .159 ± .015 | .145 ± .014 | .140 ± .015 | .141 ± .015 | .165 ± .014 | .166 ± .014 | .189 ± .015 | .189 ± .015 | .136 ± .014 |
| | .95 | .228 ± .017 | .124 ± .014 | .138 ± .015 | .125 ± .015 | .142 ± .015 | .117 ± .014 | .125 ± .014 | 0.16 ± 0.04 | .124 ± .015 | .141 ± .015 | .151 ± .015 | .134 ± .015 | .129 ± .014 | .165 ± .015 | .168 ± .014 | .171 ± .016 | .184 ± .015 | .128 ± .014 |
| | .90 | .220 ± .017 | .118 ± .014 | .124 ± .014 | .114 ± .013 | .123 ± .015 | 0.03 ± 0.04 | .121 ± .014 | .108 ± .014 | .104 ± .013 | .115 ± .014 | .156 ± .016 | .112 ± .014 | .110 ± .013 | .163 ± .015 | .171 ± .014 | .157 ± .016 | .176 ± .015 | .124 ± .015 |
| | .85 | .217 ± .017 | .101 ± .013 | .121 ± .014 | 0.01 ± 0.04 | .119 ± .014 | .092 ± .014 | .109 ± .014 | .095 ± .013 | .085 ± .013 | .107 ± .014 | .165 ± .017 | .084 ± .013 | .108 ± .013 | .161 ± .015 | .170 ± .015 | .148 ± .015 | .164 ± .016 | .107 ± .015 |
| | .80 | .214 ± .018 | .093 ± .013 | .111 ± .014 | .094 ± .014 | .109 ± .014 | .081 ± .013 | .094 ± .013 | .084 ± .013 | .087 ± .013 | .094 ± .013 | .169 ± .018 | 0.03 ± 0.03 | .169 ± .018 | 0.03 ± 0.03 | .161 ± .015 | .168 ± .015 | .158 ± .016 | .096 ± .014 |
| | .75 | .224 ± .018 | .079 ± .013 | .099 ± .014 | .080 ± .012 | .101 ± .014 | .078 ± .013 | .085 ± .014 | .079 ± .012 | .084 ± .013 | .080 ± .012 | .174 ± .019 | 0.07 ± 0.03 | .082 ± .012 | .158 ± .015 | .161 ± .015 | .136 ± .015 | .151 ± .016 | .095 ± .014 |
| | .70 | .218 ± .018 | .073 ± .012 | .097 ± .013 | .081 ± .014 | .094 ± .013 | .079 ± .012 | .079 ± .012 | .072 ± .012 | .078 ± .014 | .079 ± .012 | .178 ± .019 | .063 ± .012 | 0.03 ± 0.01 | .156 ± .015 | .132 ± .016 | .147 ± .016 | .093 ± .015 | |
| $\hat{\phi}$ | .99 | .978 ± .005 | .984 ± .005 | .992 ± .004 | .989 ± .004 | 1.000 ± .000 | .997 ± .002 | .984 ± .005 | 0.99 ± 0.03 | .979 ± .007 | .990 ± .004 | .988 ± .004 | .994 ± .003 | .989 ± .004 | .995 ± .003 | .988 ± .004 | .989 ± .005 | .994 ± .003 | 0.92 ± 0.03 |
| | .95 | .921 ± .011 | .948 ± .010 | .961 ± .008 | .941 ± .009 | .960 ± .008 | .938 ± .009 | .945 ± .009 | .929 ± .010 | .949 ± .009 | .936 ± .010 | .947 ± .009 | .975 ± .007 | .954 ± .008 | .946 ± .009 | 0.96 ± 0.08 | .968 ± .007 | .958 ± .009 | .909 ± .009 |
| | .90 | .882 ± .012 | .907 ± .012 | .927 ± .009 | .894 ± .012 | .903 ± .011 | .895 ± .011 | .909 ± .012 | .894 ± .012 | .892 ± .012 | .888 ± .012 | .898 ± .012 | .916 ± .011 | .899 ± .011 | .922 ± .010 | .902 ± .011 | .912 ± .011 | .907 ± .012 | |
| | .85 | .837 ± .014 | .861 ± .014 | .878 ± .011 | .846 ± .013 | .860 ± .012 | .837 ± .013 | .856 ± .012 | .842 ± .013 | .845 ± .013 | .839 ± .013 | .857 ± .014 | .874 ± .012 | .850 ± .012 | .867 ± .011 | .847 ± .012 | .852 ± .011 | .849 ± .013 | .804 ± .013 |
| | .80 | .791 ± .015 | .819 ± .014 | .834 ± .016 | .786 ± .015 | .799 ± .017 | .812 ± .016 | .820 ± .015 | .785 ± .016 | .787 ± .015 | .797 ± .016 | .816 ± .016 | .828 ± .014 | .814 ± .015 | .828 ± .014 | .814 ± .015 | .811 ± .016 | .811 ± .016 | |
| | .75 | .726 ± .016 | .760 ± .017 | .772 ± .016 | .767 ± .017 | .730 ± .018 | .764 ± .018 | .762 ± .017 | .734 ± .016 | .767 ± .016 | .740 ± .018 | .764 ± .018 | .754 ± .017 | .769 ± .015 | .804 ± .014 | 0.75 ± 0.04 | .805 ± .017 | .812 ± .016 | .755 ± .018 |
| | .70 | .698 ± .018 | .704 ± .018 | .691 ± .019 | .735 ± .019 | .706 ± .018 | .729 ± .019 | .685 ± .018 | .699 ± .018 | .736 ± .018 | .709 ± .019 | .729 ± .019 | .704 ± .019 | 0.65 ± 0.06 | .765 ± .017 | 0.67 ± 0.05 | .765 ± .018 | .712 ± .019 | |
| $MinCoeff$ | .99 | 1.011 ± .038 | 1.008 ± .037 | 1.008 ± .037 | 1.000 ± 0.07 | 1.000 ± 0.07 | 1.000 ± .037 | 1.008 ± .038 | 1.007 ± .037 | 1.006 ± .038 | 991 ± .037 | 1.001 ± 0.27 | 1.003 ± .038 | 1.002 ± .037 | 1.005 ± .037 | 988 ± .037 | 1.008 ± .037 | 1.001 ± 0.27 | |
| | .95 | 1.037 ± .039 | 1.008 ± .038 | 1.009 ± .038 | 1.004 ± .039 | .988 ± .037 | 1.017 ± .038 | 1.003 ± .037 | 1.019 ± .038 | .992 ± .039 | 990 ± .038 | 1.029 ± .040 | 1.004 ± 0.27 | 1.005 ± .037 | 1.002 ± .040 | .974 ± .037 | 1.029 ± .038 | 1.015 ± .039 | 1.026 ± .037 |
| | .90 | 1.032 ± .039 | .998 ± .039 | 1.018 ± .038 | 1.012 ± .038 | 0.99 ± 0.07 | 1.023 ± .039 | .998 ± .038 | 1.011 ± .038 | 1.029 ± .038 | .975 ± .039 | 1.086 ± .038 | 1.000 ± .038 | 1.016 ± .037 | 1.063 ± .040 | .975 ± .037 | 1.029 ± .039 | 1.052 ± .039 | 1.045 ± .037 |
| | .85 | 1.010 ± .040 | .987 ± .040 | 1.017 ± .039 | 0.99 ± 0.08 | .974 ± .039 | 1.013 ± .039 | .998 ± .039 | 1.019 ± .038 | .979 ± .038 | 1.084 ± .040 | 1.020 ± .038 | .999 ± .039 | 1.089 ± .042 | .967 ± .036 | 1.049 ± .040 | 1.085 ± .040 | 1.075 ± .039 | |
| | .80 | 1.002 ± .041 | .992 ± .041 | 1.019 ± .041 | .989 ± .040 | .987 ± .039 | 1.007 ± .040 | .964 ± .041 | 1.022 ± .042 | .988 ± .040 | .971 ± .039 | 1.110 ± .042 | 1.006 ± 0.40 | .998 ± .041 | 1.128 ± .045 | .959 ± .047 | 1.065 ± .041 | 1.113 ± .042 | 1.104 ± .039 |
| | .75 | .901 ± .042 | .995 ± .042 | .987 ± .042 | 0.99 ± 0.08 | .996 ± .040 | 1.005 ± .042 | .967 ± .042 | 1.012 ± .042 | .990 ± .038 | .967 ± .039 | 1.124 ± .042 | 1.008 ± .040 | .988 ± .042 | 1.155 ± .044 | .989 ± .039 | 1.073 ± .041 | 1.119 ± .042 | 1.147 ± .042 |
| | .70 | .938 ± .043 | .995 ± .044 | .965 ± .045 | 1.006 ± .042 | 1.004 ± 0.40 | 1.007 ± .042 | .958 ± .043 | .999 ± .041 | 1.055 ± .042 | .951 ± .041 | 1.118 ± .044 | 1.019 ± .043 | .984 ± .044 | 1.259 ± .047 | .910 ± .041 | 1.069 ± .043 | 1.143 ± .042 | 1.165 ± .043 |

Table B35: Results for pneumoniast: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfiNet | SELE | REG | SCross | AUCross | PlugInAUC |
|-----------------|-----|-------------|-------------|-------------|--------------------|--------------------|-------------|-------------|-------------|-------------|--------------|--------------------|--------------------|--------------------|-------------|-------------|--------------------|-------------|-------------|
| \widehat{Err} | .99 | .038 ± .006 | .033 ± .005 | .041 ± .006 | .039 ± .006 | .042 ± .006 | .033 ± .005 | .032 ± .005 | .041 ± .006 | .037 ± .006 | .045 ± .006 | .031 ± .005 | .030 ± .005 | 0.28 ± 0.05 | .040 ± .006 | .047 ± .006 | .045 ± .006 | .048 ± .006 | .033 ± .005 |
| | .95 | .035 ± .006 | .029 ± .005 | .029 ± .006 | .025 ± .005 | .035 ± .006 | .037 ± .006 | .021 ± .005 | .024 ± .005 | .022 ± .004 | .035 ± .006 | .022 ± .005 | 0.12 ± 0.03 | .019 ± .004 | .040 ± .006 | .037 ± .006 | .032 ± .006 | .031 ± .005 | .019 ± .004 |
| | .90 | .025 ± .005 | .015 ± .004 | .018 ± .005 | 0.07 ± 0.02 | 0.07 ± 0.02 | .016 ± .004 | .013 ± .004 | .012 ± .004 | .012 ± .004 | .018 ± .005 | .014 ± .004 | .014 ± .004 | .018 ± .006 | .045 ± .007 | .032 ± .005 | .038 ± .004 | .018 ± .004 | .009 ± .003 |
| | .85 | .019 ± .005 | .008 ± .003 | .011 ± .004 | .006 ± .003 | .009 ± .003 | .006 ± .002 | .008 ± .003 | .012 ± .004 | .006 ± .003 | .010 ± .003 | .005 ± .002 | 0.04 ± 0.02 | .010 ± .004 | .037 ± .007 | .045 ± .007 | .033 ± .004 | .009 ± .003 | .005 ± .002 |
| | .80 | .015 ± .004 | .005 ± .002 | .009 ± .004 | .000 ± .002 | .007 ± .003 | .005 ± .002 | .005 ± .002 | .009 ± .004 | .004 ± .002 | .002 ± .003 | 0.03 ± 0.02 | .010 ± .004 | .037 ± .007 | .045 ± .007 | .007 ± .003 | .008 ± .003 | .004 ± .002 | |
| | .75 | .009 ± .003 | .002 ± .001 | .006 ± .003 | .002 ± .001 | .008 ± .003 | .004 ± .002 | .003 ± .002 | .009 ± .003 | .002 ± .001 | .004 ± .008 | 0.01 ± 0.01 | .002 ± .002 | .010 ± .004 | .035 ± .007 | .046 ± .008 | 0.01 ± 0.01 | .006 ± .003 | .002 ± .001 |
| | .70 | .008 ± .003 | .002 ± .002 | .004 ± .002 | .002 ± .002 | .005 ± .002 | .002 ± .002 | .001 ± .001 | .004 ± .002 | .003 ± .002 | .003 ± .009 | 0.00 ± 0.00 | 0.00 ± 0.00 | .008 ± .004 | .033 ± .007 | .047 ± .008 | 0.00 ± 0.00 | .004 ± .002 | . |

DEEP NEURAL NETWORK BENCHMARKS FOR SELECTIVE CLASSIFICATION

Table B38: Results for r1: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhnGNAUC | |
|-----------------|--------------|--------------|--------------|-------------|--------------|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|--------------|-------------|--------------|--------------|--------------|-------------|
| \widehat{Err} | 99 | 499 ± 0.16 | 289 ± 0.15 | 271 ± 0.14 | 269 ± 0.14 | 295 ± 0.15 | 261 ± 0.15 | 292 ± 0.15 | 273 ± 0.14 | 256 ± 0.14 | 294 ± 0.15 | 244 ± 0.16 | 236 ± 0.15 | 206 ± 0.14 | 331 ± 0.16 | 334 ± 0.15 | 276 ± 0.16 | 274 ± 0.16 | 263 ± 0.15 | |
| | 95 | 497 ± 0.16 | 287 ± 0.15 | 274 ± 0.15 | 276 ± 0.16 | 281 ± 0.15 | 261 ± 0.15 | 292 ± 0.15 | 273 ± 0.14 | 256 ± 0.14 | 294 ± 0.15 | 244 ± 0.16 | 236 ± 0.15 | 206 ± 0.14 | 331 ± 0.16 | 334 ± 0.15 | 276 ± 0.16 | 274 ± 0.16 | 263 ± 0.15 | |
| | 90 | 499 ± 0.18 | 274 ± 0.15 | 259 ± 0.15 | 264 ± 0.16 | 276 ± 0.16 | 246 ± 0.16 | 276 ± 0.17 | 260 ± 0.15 | 257 ± 0.17 | 272 ± 0.17 | 229 ± 0.16 | 216 ± 0.16 | 229 ± 0.15 | 325 ± 0.15 | 330 ± 0.15 | 245 ± 0.16 | 222 ± 0.17 | 244 ± 0.16 | |
| | 85 | 499 ± 0.19 | 270 ± 0.16 | 248 ± 0.16 | 269 ± 0.15 | 280 ± 0.18 | 229 ± 0.16 | 267 ± 0.17 | 248 ± 0.15 | 250 ± 0.15 | 265 ± 0.16 | 234 ± 0.16 | 209 ± 0.16 | 246 ± 0.15 | 325 ± 0.16 | 334 ± 0.16 | 236 ± 0.16 | 224 ± 0.17 | 233 ± 0.16 | |
| | 80 | 506 ± 0.20 | 266 ± 0.17 | 257 ± 0.17 | 231 ± 0.17 | 253 ± 0.16 | 217 ± 0.16 | 258 ± 0.17 | 235 ± 0.16 | 229 ± 0.17 | 238 ± 0.17 | 204 ± 0.17 | 233 ± 0.15 | 223 ± 0.16 | 321 ± 0.16 | 315 ± 0.16 | 221 ± 0.16 | 221 ± 0.18 | 239 ± 0.17 | |
| | 75 | 513 ± 0.20 | 262 ± 0.16 | 213 ± 0.16 | 245 ± 0.16 | 252 ± 0.18 | 198 ± 0.16 | 249 ± 0.18 | 217 ± 0.16 | 224 ± 0.16 | 232 ± 0.17 | 236 ± 0.17 | 192 ± 0.16 | 222 ± 0.16 | 319 ± 0.16 | 302 ± 0.16 | 200 ± 0.17 | 211 ± 0.18 | 223 ± 0.17 | |
| | 70 | 510 ± 0.21 | 256 ± 0.17 | 207 ± 0.17 | 224 ± 0.18 | 241 ± 0.18 | 195 ± 0.17 | 231 ± 0.19 | 202 ± 0.17 | 219 ± 0.17 | 233 ± 0.17 | 227 ± 0.17 | 175 ± 0.16 | 215 ± 0.15 | 319 ± 0.17 | 293 ± 0.16 | 188 ± 0.18 | 193 ± 0.17 | 205 ± 0.17 | |
| | $\hat{\phi}$ | 99 | 992 ± 0.03 | 991 ± 0.03 | 986 ± 0.03 | 992 ± 0.03 | 982 ± 0.04 | 984 ± 0.03 | 998 ± 0.02 | 992 ± 0.03 | 994 ± 0.03 | 984 ± 0.04 | 990 ± 0.03 | 971 ± 0.05 | 980 ± 0.05 | 986 ± 0.04 | 988 ± 0.04 | 979 ± 0.05 | 987 ± 0.03 | 994 ± 0.02 |
| | | 95 | 974 ± 0.02 | 942 ± 0.07 | 939 ± 0.08 | 942 ± 0.08 | 950 ± 0.07 | 943 ± 0.07 | 959 ± 0.06 | 950 ± 0.07 | 973 ± 0.05 | 943 ± 0.07 | 967 ± 0.05 | 927 ± 0.08 | 944 ± 0.08 | 955 ± 0.07 | 968 ± 0.05 | 965 ± 0.08 | 919 ± 0.09 | 934 ± 0.08 |
| | | 90 | 969 ± 0.009 | 901 ± 0.009 | 896 ± 0.010 | 898 ± 0.010 | 902 ± 0.010 | 902 ± 0.008 | 897 ± 0.010 | 897 ± 0.010 | 879 ± 0.010 | 881 ± 0.010 | 914 ± 0.009 | 885 ± 0.009 | 878 ± 0.010 | 926 ± 0.009 | 928 ± 0.007 | 858 ± 0.012 | 844 ± 0.012 | 878 ± 0.011 |
| 85 | | 871 ± 0.10 | 843 ± 0.11 | 842 ± 0.12 | 858 ± 0.14 | 845 ± 0.11 | 846 ± 0.11 | 853 ± 0.12 | 839 ± 0.11 | 831 ± 0.13 | 847 ± 0.10 | 854 ± 0.12 | 848 ± 0.11 | 840 ± 0.12 | 880 ± 0.11 | 869 ± 0.10 | 765 ± 0.13 | 764 ± 0.13 | 821 ± 0.12 | |
| 80 | | 822 ± 0.12 | 812 ± 0.12 | 757 ± 0.11 | 814 ± 0.11 | 802 ± 0.14 | 794 ± 0.13 | 803 ± 0.12 | 777 ± 0.14 | 815 ± 0.12 | 801 ± 0.13 | 804 ± 0.14 | 786 ± 0.13 | 788 ± 0.14 | 834 ± 0.11 | 830 ± 0.11 | 703 ± 0.14 | 704 ± 0.14 | 800 ± 0.13 | |
| 75 | | 785 ± 0.13 | 779 ± 0.14 | 716 ± 0.15 | 781 ± 0.12 | 753 ± 0.14 | 739 ± 0.14 | 752 ± 0.13 | 730 ± 0.15 | 744 ± 0.15 | 739 ± 0.14 | 770 ± 0.13 | 739 ± 0.14 | 746 ± 0.13 | 753 ± 0.13 | 757 ± 0.13 | 647 ± 0.15 | 667 ± 0.16 | 759 ± 0.13 | |
| 70 | | 728 ± 0.15 | 723 ± 0.15 | 670 ± 0.16 | 678 ± 0.16 | 718 ± 0.15 | 701 ± 0.14 | 704 ± 0.14 | 687 ± 0.16 | 700 ± 0.15 | 705 ± 0.15 | 724 ± 0.14 | 691 ± 0.15 | 696 ± 0.15 | 727 ± 0.13 | 713 ± 0.14 | 608 ± 0.17 | 608 ± 0.17 | 713 ± 0.15 | |
| $MinCoeff$ | | 99 | 999 ± 0.02 | 1004 ± 0.02 | 999 ± 0.03 | 997 ± 0.03 | 999 ± 0.03 | 1002 ± 0.03 | 1003 ± 0.02 | 1001 ± 0.03 | 1001 ± 0.02 | 998 ± 0.02 | 1009 ± 0.03 | 1005 ± 0.03 | 1007 ± 0.03 | 999 ± 0.03 | 1004 ± 0.02 | 999 ± 0.02 | 1004 ± 0.02 | 1005 ± 0.03 |
| | | 95 | 994 ± 0.03 | 998 ± 0.03 | 996 ± 0.03 | 1008 ± 0.04 | 992 ± 0.04 | 1003 ± 0.02 | 1003 ± 0.03 | 1007 ± 0.04 | 1013 ± 0.03 | 1002 ± 0.03 | 1018 ± 0.04 | 1008 ± 0.04 | 999 ± 0.04 | 997 ± 0.04 | 995 ± 0.02 | 989 ± 0.03 | 1022 ± 0.04 | 1017 ± 0.04 |
| | | 90 | 998 ± 0.036 | 993 ± 0.03 | 992 ± 0.03 | 1011 ± 0.03 | 992 ± 0.05 | 1000 ± 0.03 | 1008 ± 0.05 | 1000 ± 0.05 | 1019 ± 0.04 | 1021 ± 0.04 | 1030 ± 0.06 | 1020 ± 0.04 | 995 ± 0.06 | 998 ± 0.05 | 999 ± 0.03 | 987 ± 0.03 | 1040 ± 0.06 | 1031 ± 0.05 |
| | 85 | 998 ± 0.038 | 993 ± 0.033 | 990 ± 0.034 | 1009 ± 0.036 | 987 ± 0.036 | 1009 ± 0.036 | 992 ± 0.036 | 994 ± 0.036 | 1021 ± 0.03 | 1014 ± 0.034 | 1031 ± 0.05 | 1019 ± 0.05 | 993 ± 0.037 | 1006 ± 0.035 | 994 ± 0.037 | 1006 ± 0.035 | 1048 ± 0.039 | 1040 ± 0.035 | |
| | 80 | 1012 ± 0.039 | 999 ± 0.033 | 984 ± 0.038 | 1014 ± 0.036 | 990 ± 0.038 | 1016 ± 0.037 | 992 ± 0.037 | 1002 ± 0.038 | 993 ± 0.035 | 989 ± 0.035 | 1053 ± 0.036 | 1040 ± 0.035 | 975 ± 0.038 | 1021 ± 0.035 | 983 ± 0.037 | 1001 ± 0.040 | 1023 ± 0.040 | 1038 ± 0.035 | |
| | 75 | 1008 ± 0.040 | 1003 ± 0.034 | 952 ± 0.032 | 1005 ± 0.035 | 996 ± 0.038 | 1028 ± 0.039 | 1000 ± 0.038 | 983 ± 0.038 | 1022 ± 0.037 | 1016 ± 0.037 | 1062 ± 0.037 | 1045 ± 0.035 | 981 ± 0.038 | 1012 ± 0.037 | 974 ± 0.038 | 927 ± 0.040 | 1036 ± 0.040 | 1049 ± 0.036 | |
| | 70 | 1019 ± 0.042 | 1008 ± 0.036 | 979 ± 0.039 | 1033 ± 0.039 | 990 ± 0.039 | 1027 ± 0.039 | 986 ± 0.038 | 985 ± 0.037 | 985 ± 0.038 | 971 ± 0.036 | 1073 ± 0.038 | 1045 ± 0.037 | 992 ± 0.040 | 1007 ± 0.037 | 971 ± 0.040 | 985 ± 0.042 | 1062 ± 0.040 | 1062 ± 0.038 | |

Table B39: Results for stanfordcars: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | |
|-----------------|--------------|------------|------------|------------|------------|------------|------------|------------|------------|------------|--------------|------------|------------|------------|------------|------------|------------|------------|
| \widehat{Err} | 99 | 236 ± 0.08 | 178 ± 0.07 | 170 ± 0.06 | 318 ± 0.09 | 469 ± 0.08 | 167 ± 0.06 | 173 ± 0.07 | 165 ± 0.06 | 317 ± 0.08 | 469 ± 0.08 | 107 ± 0.05 | 109 ± 0.06 | 230 ± 0.08 | 429 ± 0.09 | 437 ± 0.10 | 102 ± 0.07 | |
| | 95 | 223 ± 0.08 | 165 ± 0.07 | 150 ± 0.06 | 285 ± 0.07 | 442 ± 0.09 | 145 ± 0.06 | 160 ± 0.07 | 143 ± 0.06 | 283 ± 0.08 | 430 ± 0.09 | 99 ± 0.05 | 106 ± 0.05 | 224 ± 0.08 | 416 ± 0.09 | 433 ± 0.09 | 141 ± 0.07 | |
| | 90 | 216 ± 0.08 | 146 ± 0.07 | 123 ± 0.06 | 281 ± 0.08 | 433 ± 0.08 | 110 ± 0.05 | 131 ± 0.06 | 115 ± 0.06 | 245 ± 0.07 | 411 ± 0.08 | 97 ± 0.05 | 105 ± 0.05 | 164 ± 0.04 | 211 ± 0.08 | 405 ± 0.09 | 431 ± 0.10 | 123 ± 0.06 |
| | 85 | 190 ± 0.08 | 127 ± 0.06 | 109 ± 0.06 | 276 ± 0.09 | 378 ± 0.10 | 96 ± 0.05 | 114 ± 0.06 | 97 ± 0.05 | 224 ± 0.08 | 349 ± 0.10 | 97 ± 0.04 | 104 ± 0.04 | 199 ± 0.08 | 395 ± 0.10 | 433 ± 0.10 | 100 ± 0.06 | |
| | 80 | 176 ± 0.08 | 111 ± 0.06 | 98 ± 0.06 | 274 ± 0.08 | 383 ± 0.10 | 87 ± 0.05 | 109 ± 0.06 | 87 ± 0.05 | 208 ± 0.08 | 372 ± 0.10 | 90 ± 0.05 | 103 ± 0.04 | 182 ± 0.08 | 382 ± 0.10 | 437 ± 0.10 | 95 ± 0.05 | |
| | 75 | 157 ± 0.08 | 109 ± 0.06 | 89 ± 0.06 | 268 ± 0.09 | 377 ± 0.08 | 98 ± 0.04 | 101 ± 0.06 | 90 ± 0.04 | 174 ± 0.08 | 340 ± 0.09 | 95 ± 0.05 | 102 ± 0.04 | 173 ± 0.08 | 373 ± 0.10 | 435 ± 0.10 | 97 ± 0.05 | |
| | 70 | 148 ± 0.08 | 101 ± 0.06 | 87 ± 0.05 | 256 ± 0.09 | 316 ± 0.09 | 107 ± 0.04 | 108 ± 0.05 | 106 ± 0.04 | 166 ± 0.08 | 327 ± 0.10 | 95 ± 0.04 | 102 ± 0.05 | 164 ± 0.08 | 362 ± 0.11 | 438 ± 0.11 | 95 ± 0.05 | |
| | $\hat{\phi}$ | 99 | 988 ± 0.02 | 994 ± 0.01 | 994 ± 0.01 | 984 ± 0.02 | 989 ± 0.02 | 989 ± 0.02 | 987 ± 0.02 | 988 ± 0.02 | 983 ± 0.01 | 991 ± 0.02 | 990 ± 0.02 | 987 ± 0.02 | 986 ± 0.02 | 992 ± 0.02 | 987 ± 0.02 | 993 ± 0.01 |
| | | 95 | 958 ± 0.03 | 966 ± 0.03 | 955 ± 0.04 | 958 ± 0.04 | 948 ± 0.03 | 954 ± 0.03 | 956 ± 0.03 | 952 ± 0.03 | 954 ± 0.03 | 949 ± 0.03 | 958 ± 0.04 | 954 ± 0.04 | 956 ± 0.04 | 954 ± 0.04 | 952 ± 0.04 | 960 ± 0.03 |
| | | 90 | 910 ± 0.05 | 925 ± 0.04 | 907 ± 0.05 | 886 ± 0.06 | 900 ± 0.05 | 896 ± 0.05 | 917 ± 0.04 | 914 ± 0.04 | 911 ± 0.04 | 899 ± 0.05 | 903 ± 0.05 | 901 ± 0.05 | 919 ± 0.05 | 911 ± 0.05 | 905 ± 0.06 | 922 ± 0.04 |
| 85 | | 868 ± 0.06 | 875 ± 0.05 | 861 ± 0.06 | 840 ± 0.07 | 853 ± 0.06 | 844 ± 0.05 | 866 ± 0.06 | 857 ± 0.05 | 855 ± 0.06 | 843 ± 0.06 | 855 ± 0.06 | 852 ± 0.06 | 866 ± 0.05 | 862 ± 0.06 | 867 ± 0.06 | 875 ± 0.05 | |
| 80 | | 815 ± 0.06 | 833 ± 0.06 | 814 ± 0.06 | 804 ± 0.07 | 802 ± 0.07 | 792 ± 0.06 | 817 ± 0.07 | 804 ± 0.06 | 815 ± 0.06 | 805 ± 0.07 | 804 ± 0.06 | 806 ± 0.07 | 813 ± 0.06 | 818 ± 0.07 | 810 ± 0.07 | 833 ± 0.06 | |
| 75 | | 761 ± 0.07 | 789 ± 0.07 | 775 ± 0.07 | 761 ± 0.08 | 757 ± 0.07 | 748 ± 0.07 | 776 ± 0.08 | 755 ± 0.07 | 764 ± 0.08 | 761 ± 0.08 | 756 ± 0.07 | 751 ± 0.07 | 768 ± 0.07 | 766 ± 0.06 | 756 ± 0.08 | 790 ± 0.07 | |
| 70 | | 715 ± 0.07 | 733 ± 0.08 | 716 ± 0.07 | 712 ± 0.07 | 713 ± 0.07 | 708 ± 0.07 | 730 ± 0.08 | 707 ± 0.07 | 716 ± 0.08 | 715 ± 0.08 | 706 ± 0.07 | 706 ± 0.07 | 727 ± 0.07 | 718 ± 0.07 | 715 ± 0.08 | 748 ± 0.07 | |

Table B40: Results for SVHN: $mean \pm std$ for \widehat{Err} and empirical coverage $\hat{\phi}$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross |
|-----------------|----|------------|------------|------------|------------|------------|------------|------------|------------|------------|--------------|------------|------------|------------|------------|------------|------------|
| \widehat{Err} | 99 | 038 ± 0.02 | 030 ± 0.01 | 037 ± 0.01 | 030 ± 0.01 | 043 ± 0.02 | 039 ± 0.01 | 030 ± 0.01 | 036 ± 0.01 | 031 ± 0.01 | 043 ± 0.02 | 036 ± 0.01 | 034 ± 0.01 | 036 ± 0.01 | 050 ± 0.02 | 050 ± 0.02 | 038 ± 0.02 |
| | 95 | 022 ± 0.01 | 017 ± 0.01 | 020 ± 0.01 | 016 ± 0.01 | 026 ± 0.01 | 021 ± 0.01 | 017 ± 0.01 | 020 ± 0.01 | 016 ± 0.01 | 024 ± 0.01 | 021 ± 0.01 | 018 ± 0.01 | 019 ± 0.01 | 049 ± 0.02 | 050 ± 0.02 | 021 ± 0.01 |
| | 90 | 013 ± 0.01 | 008 ± 0.01 | 011 ± 0.01 | 009 ± 0.01 | 014 ± 0.01 | 011 ± 0.01 | 008 ± 0.01 | 010 ± 0.01 | 009 ± 0.01 | 013 ± 0.01 | 011 ± 0.01 | 010 ± 0.01 | 009 ± 0.01 | 048 ± 0.02 | 050 ± 0.02 | 011 ± 0.01 |
| | 85 | 012 ± 0.01 | 009 ± 0.01 | 009 ± 0.01 | 009 ± 0.01 | 009 ± 0.01 | 011 ± 0.01 | 007 ± 0.01 | 008 ± 0.01 | 007 ± 0.01 | 007 ± 0.01 | 007 ± 0.01 | 007 ± 0.01 | 007 ± 0.01 | 045 ± 0.02 | 045 ± 0.02 | 008 ± 0.02 |
| | 80 | 011 ± 0.01 | 005 ± 0.01 | 006 ± 0.01 | 005 ± 0.01 | 008 ± 0.01 | 006 ± 0.01 | 005 ± 0.01 | 006 ± 0.01 | 005 ± 0.01 | 007 ± 0.01 | 006 ± 0.01 | 005 ± 0.01 | 006 ± 0.01 | 042 ± 0.02 | 050 ± 0.02 | 006 ± 0.01 |
| | 75 | 011 ± 0.01 | 004 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 006 ± 0.01 | 037 ± 0.02 | 051 ± 0.02 | 005 ± 0.01 |
| | 70 | 010 ± 0.01 | 004 ± 0.01 | 005 ± 0.01 | 005 ± 0.01 | 004 ± 0.01 | 004 ± | | | | | | | | | | |

Table B43: Results for upselling: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhgInAUC | |
|-----------------|----|-------------------|----------------------|-------------------|----------------------|----------------------|---------------|--------------------|-------------------|----------------------|---------------|-------------------|-------------------|--------------------|----------------------|-------------------|---------------|----------------------|---------------|---------------|
| \widehat{Err} | 99 | 277 ± 0.04 | 168 ± 0.03 | 171 ± 0.04 | 164 ± 0.03 | 178 ± 0.04 | 186 ± 0.05 | 169 ± 0.04 | 171 ± 0.04 | 161 ± 0.03 | 179 ± 0.03 | 189 ± 0.03 | 186 ± 0.03 | 177 ± 0.02 | 185 ± 0.04 | 185 ± 0.03 | 177 ± 0.02 | 185 ± 0.02 | 191 ± 0.05 | |
| | 95 | 278 ± 0.04 | 169 ± 0.03 | 163 ± 0.03 | 162 ± 0.03 | 171 ± 0.04 | 174 ± 0.05 | 163 ± 0.03 | 159 ± 0.03 | 152 ± 0.03 | 171 ± 0.03 | 185 ± 0.04 | 168 ± 0.03 | 166 ± 0.02 | 179 ± 0.04 | 178 ± 0.03 | 177 ± 0.02 | 185 ± 0.02 | 183 ± 0.05 | |
| | 90 | 276 ± 0.05 | 150 ± 0.04 | 150 ± 0.03 | 145 ± 0.02 | 155 ± 0.03 | 163 ± 0.04 | 142 ± 0.03 | 142 ± 0.03 | 156 ± 0.03 | 208 ± 0.04 | 156 ± 0.03 | 208 ± 0.04 | 156 ± 0.03 | 148 ± 0.02 | 164 ± 0.04 | 172 ± 0.04 | 177 ± 0.02 | 185 ± 0.02 | 164 ± 0.05 |
| | 85 | 275 ± 0.05 | 136 ± 0.03 | 134 ± 0.03 | 122 ± 0.02 | 139 ± 0.03 | 145 ± 0.04 | 126 ± 0.03 | 129 ± 0.02 | 128 ± 0.03 | 137 ± 0.03 | 221 ± 0.05 | 142 ± 0.03 | 138 ± 0.02 | 162 ± 0.04 | 168 ± 0.04 | 177 ± 0.02 | 185 ± 0.02 | 157 ± 0.05 | |
| | 80 | 291 ± 0.06 | 104 ± 0.02 | 121 ± 0.03 | 115 ± 0.02 | 121 ± 0.02 | 130 ± 0.03 | 107 ± 0.02 | 116 ± 0.02 | 110 ± 0.03 | 124 ± 0.02 | 230 ± 0.06 | 118 ± 0.01 | 114 ± 0.02 | 144 ± 0.04 | 168 ± 0.05 | 176 ± 0.02 | 185 ± 0.02 | 147 ± 0.04 | |
| | 75 | 297 ± 0.06 | 094 ± 0.01 | 110 ± 0.02 | 104 ± 0.02 | 102 ± 0.01 | 106 ± 0.02 | 092 ± 0.02 | 103 ± 0.02 | 101 ± 0.02 | 105 ± 0.02 | 244 ± 0.06 | 095 ± 0.01 | 098 ± 0.02 | 128 ± 0.04 | 164 ± 0.04 | 171 ± 0.02 | 185 ± 0.02 | 123 ± 0.03 | |
| | 70 | 304 ± 0.07 | 078 ± 0.01 | 081 ± 0.00 | 075 ± 0.00 | 080 ± 0.00 | 080 ± 0.00 | 080 ± 0.01 | 080 ± 0.01 | 077 ± 0.01 | 077 ± 0.01 | 257 ± 0.07 | 076 ± 0.01 | 074 ± 0.00 | 120 ± 0.03 | 158 ± 0.04 | 171 ± 0.02 | 185 ± 0.02 | 108 ± 0.03 | |
| $\hat{\phi}$ | 99 | 994 ± 0.02 | 988 ± 0.04 | 985 ± 0.04 | 993 ± 0.03 | 984 ± 0.04 | 989 ± 0.04 | 988 ± 0.04 | 995 ± 0.02 | 984 ± 0.04 | 994 ± 0.02 | 992 ± 0.03 | 990 ± 0.03 | 993 ± 0.03 | 990 ± 0.03 | 991 ± 0.03 | 978 ± 0.05 | 998 ± 0.01 | 996 ± 0.02 | |
| | 95 | 948 ± 0.07 | 950 ± 0.07 | 951 ± 0.07 | 964 ± 0.06 | 950 ± 0.08 | 951 ± 0.08 | 965 ± 0.06 | 953 ± 0.07 | 940 ± 0.09 | 944 ± 0.07 | 961 ± 0.06 | 941 ± 0.07 | 953 ± 0.06 | 962 ± 0.06 | 939 ± 0.07 | 975 ± 0.05 | 998 ± 0.01 | 953 ± 0.06 | |
| | 90 | 888 ± 0.10 | 906 ± 0.10 | 883 ± 0.11 | 915 ± 0.10 | 918 ± 0.09 | 923 ± 0.10 | 884 ± 0.12 | 887 ± 0.11 | 902 ± 0.11 | 912 ± 0.10 | 900 ± 0.10 | 895 ± 0.10 | 913 ± 0.10 | 907 ± 0.08 | 881 ± 0.09 | 978 ± 0.05 | 998 ± 0.01 | 890 ± 0.10 | |
| | 85 | 833 ± 0.12 | 860 ± 0.12 | 842 ± 0.12 | 841 ± 0.13 | 863 ± 0.12 | 879 ± 0.11 | 841 ± 0.13 | 849 ± 0.13 | 851 ± 0.13 | 848 ± 0.12 | 847 ± 0.11 | 858 ± 0.12 | 867 ± 0.12 | 881 ± 0.09 | 852 ± 0.10 | 978 ± 0.05 | 998 ± 0.01 | 872 ± 0.11 | |
| | 80 | 759 ± 0.14 | 808 ± 0.14 | 806 ± 0.13 | 805 ± 0.13 | 801 ± 0.14 | 833 ± 0.12 | 805 ± 0.14 | 807 ± 0.14 | 795 ± 0.15 | 816 ± 0.13 | 810 ± 0.12 | 809 ± 0.13 | 822 ± 0.14 | 823 ± 0.11 | 798 ± 0.12 | 977 ± 0.05 | 998 ± 0.01 | 825 ± 0.12 | |
| | 75 | 718 ± 0.15 | 764 ± 0.14 | 758 ± 0.14 | 772 ± 0.14 | 758 ± 0.14 | 786 ± 0.14 | 758 ± 0.15 | 768 ± 0.14 | 769 ± 0.15 | 760 ± 0.14 | 760 ± 0.12 | 768 ± 0.15 | 772 ± 0.15 | 762 ± 0.14 | 754 ± 0.13 | 964 ± 0.06 | 998 ± 0.01 | 772 ± 0.13 | |
| | 70 | 677 ± 0.15 | 724 ± 0.15 | 702 ± 0.15 | 722 ± 0.15 | 706 ± 0.16 | 725 ± 0.15 | 721 ± 0.16 | 739 ± 0.15 | 722 ± 0.16 | 724 ± 0.15 | 722 ± 0.14 | 709 ± 0.15 | 723 ± 0.15 | 730 ± 0.14 | 717 ± 0.12 | 964 ± 0.06 | 998 ± 0.01 | 721 ± 0.14 | |
| $MinCoeff$ | 99 | 1.000 ± 0.031 | 1.003 ± 0.031 | 1.004 ± 0.031 | 1.002 ± 0.031 | 1.001 ± 0.031 | 1.001 ± 0.030 | 1.005 ± 0.030 | 1.002 ± 0.030 | 1.000 ± 0.031 | 999 ± 0.030 | 1.009 ± 0.031 | 1.001 ± 0.031 | 1.003 ± 0.030 | 1.001 ± 0.031 | 999 ± 0.031 | 1.003 ± 0.031 | 1.001 ± 0.031 | 1.001 ± 0.031 | |
| | 95 | 977 ± 0.032 | 1.000 ± 0.031 | 1.004 ± 0.033 | 1.001 ± 0.031 | 1.001 ± 0.031 | 995 ± 0.031 | 1.005 ± 0.031 | 1.016 ± 0.030 | 1.004 ± 0.031 | 1.000 ± 0.031 | 985 ± 0.032 | 1.011 ± 0.031 | 997 ± 0.031 | 1.005 ± 0.031 | 1.011 ± 0.032 | 1.003 ± 0.031 | 1.001 ± 0.031 | 1.013 ± 0.031 | |
| | 90 | 959 ± 0.033 | 998 ± 0.031 | 1.017 ± 0.034 | 994 ± 0.032 | 1.004 ± 0.031 | 997 ± 0.032 | 993 ± 0.032 | 1.020 ± 0.031 | 1.001 ± 0.032 | 1.002 ± 0.031 | 997 ± 0.031 | 1.012 ± 0.031 | 996 ± 0.031 | 1.027 ± 0.032 | 1.007 ± 0.032 | 1.003 ± 0.031 | 1.001 ± 0.031 | 1.021 ± 0.032 | |
| | 85 | 921 ± 0.034 | 998 ± 0.031 | 1.018 ± 0.035 | 991 ± 0.033 | 1.005 ± 0.033 | 1.004 ± 0.033 | 996 ± 0.032 | 1.024 ± 0.031 | 1.024 ± 0.031 | 992 ± 0.032 | 1.018 ± 0.032 | 985 ± 0.035 | 1.006 ± 0.031 | 1.005 ± 0.030 | 1.035 ± 0.032 | 998 ± 0.033 | 1.003 ± 0.031 | 1.001 ± 0.031 | 1.022 ± 0.032 |
| | 80 | 849 ± 0.034 | 1.008 ± 0.032 | 1.022 ± 0.034 | 1.007 ± 0.032 | 1.011 ± 0.032 | 1.007 ± 0.033 | 997 ± 0.032 | 1.026 ± 0.032 | 1.009 ± 0.033 | 1.004 ± 0.033 | 829 ± 0.035 | 1.012 ± 0.032 | 997 ± 0.032 | 1.048 ± 0.034 | 984 ± 0.033 | 1.002 ± 0.031 | 1.001 ± 0.031 | 1.041 ± 0.033 | |
| | 75 | 799 ± 0.034 | 1.018 ± 0.034 | 1.021 ± 0.036 | 1.019 ± 0.033 | 1.013 ± 0.033 | 1.019 ± 0.034 | 1.013 ± 0.033 | 1.015 ± 0.033 | 1.016 ± 0.034 | 1.016 ± 0.033 | 781 ± 0.035 | 1.012 ± 0.033 | 1.015 ± 0.033 | 1.077 ± 0.034 | 976 ± 0.034 | 1.000 ± 0.032 | 1.001 ± 0.031 | 1.050 ± 0.035 | |
| | 70 | 746 ± 0.036 | 1.013 ± 0.035 | 1.039 ± 0.037 | 1.026 ± 0.034 | 1.042 ± 0.035 | 1.029 ± 0.035 | 1.020 ± 0.035 | 1.025 ± 0.034 | 1.017 ± 0.035 | 1.022 ± 0.033 | 781 ± 0.036 | 1.023 ± 0.036 | 1.024 ± 0.034 | 1.095 ± 0.036 | 965 ± 0.035 | 1.000 ± 0.032 | 1.001 ± 0.031 | 1.066 ± 0.037 | |

Table B44: Results for waterbirds: $mean \pm std$ for \widehat{Err} , empirical coverage $\hat{\phi}$, and $MinCoeff$.

| Metric | c | DG | SAT | SAT+EM | SelNet | SelNet+EM | SR | SAT+SR | SAT+EM+SR | SelNet+SR | SelNet+EM+SR | ENS | ENS+SR | ConfidNet | SELE | REG | SCross | AUCross | PhgInAUC |
|-----------------|----|-------------------|-------------------|-------------|-------------------|-------------------|-------------------|-------------|-------------|--------------------|-------------------|--------------------|-------------------|-------------------|-------------|----------------------|--------------------|---------------|---------------|
| \widehat{Err} | 99 | 143 ± 0.08 | 093 ± 0.06 | 109 ± 0.07 | 114 ± 0.07 | 120 ± 0.07 | 094 ± 0.06 | 094 ± 0.06 | 110 ± 0.07 | 115 ± 0.07 | 117 ± 0.07 | 083 ± 0.06 | 083 ± 0.06 | 101 ± 0.06 | 139 ± 0.07 | 144 ± 0.07 | 102 ± 0.06 | 108 ± 0.06 | 099 ± 0.07 |
| | 95 | 142 ± 0.08 | 078 ± 0.06 | 094 ± 0.07 | 109 ± 0.06 | 113 ± 0.07 | 078 ± 0.06 | 080 ± 0.06 | 090 ± 0.07 | 102 ± 0.07 | 119 ± 0.07 | 075 ± 0.06 | 084 ± 0.06 | 093 ± 0.06 | 137 ± 0.07 | 141 ± 0.07 | 084 ± 0.05 | 110 ± 0.06 | 097 ± 0.07 |
| | 90 | 134 ± 0.08 | 064 ± 0.05 | 085 ± 0.07 | 101 ± 0.07 | 087 ± 0.07 | 068 ± 0.06 | 062 ± 0.05 | 073 ± 0.06 | 095 ± 0.07 | 091 ± 0.07 | 063 ± 0.05 | 049 ± 0.05 | 081 ± 0.06 | 133 ± 0.07 | 141 ± 0.07 | 065 ± 0.05 | 115 ± 0.06 | 098 ± 0.07 |
| | 85 | 112 ± 0.08 | 051 ± 0.05 | 073 ± 0.06 | 086 ± 0.06 | 083 ± 0.06 | 056 ± 0.05 | 049 ± 0.05 | 064 ± 0.06 | 090 ± 0.06 | 076 ± 0.06 | 052 ± 0.05 | 039 ± 0.05 | 064 ± 0.05 | 127 ± 0.07 | 139 ± 0.07 | 048 ± 0.05 | 119 ± 0.07 | 096 ± 0.07 |
| | 80 | 089 ± 0.07 | 043 ± 0.05 | 056 ± 0.05 | 088 ± 0.07 | 095 ± 0.07 | 047 ± 0.06 | 040 ± 0.05 | 053 ± 0.06 | 069 ± 0.06 | 063 ± 0.06 | 041 ± 0.05 | 031 ± 0.04 | 055 ± 0.05 | 126 ± 0.08 | 139 ± 0.08 | 035 ± 0.04 | 127 ± 0.07 | 094 ± 0.07 |
| | 75 | 072 ± 0.07 | 033 ± 0.04 | 044 ± 0.05 | 082 ± 0.06 | 080 ± 0.07 | 040 ± 0.05 | 033 ± 0.04 | 041 ± 0.05 | 065 ± 0.05 | 055 ± 0.06 | 034 ± 0.05 | 023 ± 0.04 | 049 ± 0.05 | 124 ± 0.08 | 141 ± 0.08 | 029 ± 0.04 | 130 ± 0.08 | 088 ± 0.07 |
| | 70 | 058 ± 0.06 | 027 ± 0.04 | 042 ± 0.05 | 079 ± 0.07 | 083 ± 0.06 | 034 ± 0.05 | 027 ± 0.04 | 040 ± 0.05 | 075 ± 0.06 | 049 ± 0.05 | 028 ± 0.04 | 018 ± 0.03 | 038 ± 0.04 | 119 ± 0.08 | 141 ± 0.08 | 017 ± 0.04 | 135 ± 0.08 | 083 ± 0.06 |
| $\hat{\phi}$ | 99 | 989 ± 0.02 | 991 ± 0.02 | 986 ± 0.03 | 981 ± 0.03 | 990 ± 0.02 | 991 ± 0.02 | 994 ± 0.01 | 987 ± 0.03 | 989 ± 0.02 | 985 ± 0.02 | 992 ± 0.02 | 990 ± 0.02 | 990 ± 0.02 | 973 ± 0.04 | 965 ± 0.01 | 988 ± 0.02 | 987 ± 0.02 | 994 ± 0.02 |
| | 95 | 951 ± 0.05 | 946 ± 0.05 | 936 ± 0.06 | 938 ± 0.05 | 943 ± 0.05 | 943 ± 0.05 | 959 ± 0.05 | 936 ± 0.06 | 927 ± 0.06 | 943 ± 0.05 | 949 ± 0.05 | 950 ± 0.04 | 945 ± 0.05 | 937 ± 0.05 | 956 ± 0.05 | 932 ± 0.05 | 939 ± 0.05 | 940 ± 0.06 |
| | 90 | 904 ± 0.06 | 897 ± 0.07 | 884 ± 0.08 | 899 ± 0.07 | 913 ± 0.06 | 906 ± 0.06 | 904 ± 0.07 | 879 ± 0.07 | 889 ± 0.06 | 912 ± 0.07 | 901 ± 0.06 | 900 ± 0.07 | 895 ± 0.07 | 878 ± 0.07 | 896 ± 0.06 | 857 ± 0.07 | 886 ± 0.07 | 894 ± 0.07 |
| | 85 | 849 ± 0.07 | 841 ± 0.09 | 845 ± 0.09 | 842 ± 0.09 | 835 ± 0.09 | 855 ± 0.08 | 844 ± 0.08 | 826 ± 0.08 | 842 ± 0.09 | 831 ± 0.09 | 856 ± 0.07 | 849 ± 0.09 | 845 ± 0.08 | 832 ± 0.08 | 834 ± 0.07 | 773 ± 0.08 | 821 ± 0.07 | 847 ± 0.08 |
| | 80 | 801 ± 0.09 | 785 ± 0.10 | 776 ± 0.09 | 799 ± 0.09 | 803 ± 0.09 | 804 ± 0.09 | 775 ± 0.10 | 779 ± 0.09 | 790 ± 0.09 | 801 ± 0.09 | 800 ± 0.09 | 808 ± 0.08 | 808 ± 0.08 | 761 ± 0.09 | 807 ± 0.08 | 695 ± 0.09 | 752 ± 0.09 | 807 ± 0.08 |
| | 75 | 749 ± 0.10 | 716 ± 0.10 | 723 ± 0.09 | 756 ± 0.10 | 757 ± 0.09 | 754 ± 0.10 | 717 ± 0.10 | 729 ± 0.10 | 745 ± 0.10 | 733 ± 0.10 | 746 ± 0.09 | 750 ± 0.10 | 756 ± 0.10 | 710 ± 0.10 | 739 ± 0.09 | 640 ± 0.09 | 685 ± 0.09 | 760 ± 0.09 |
| | 70 | 705 ± 0.10 | 677 ± 0.11 | 684 ± 0.10 | 693 ± 0.11 | 710 ± 0.10 | 699 ± 0.11 | 671 ± 0.11 | 689 ± 0.10 | 697 ± 0.09 | 717 ± 0.10 | 698 ± 0.09 | 699 ± 0.11 | 690 ± 0.10 | 664 ± 0.10 | 707 ± 0.09 | 573 ± 0.09 | 626 ± 0.10 | 701 ± 0.10 |
| $MinCoeff$ | 99 | 988 ± 0.037 | 987 ± 0.039 | 987 ± 0.039 | 988 ± 0.037 | 980 ± 0.037 | 981 ± 0.038 | 993 ± 0.038 | 982 ± 0.037 | 997 ± 0.038 | 989 ± 0.037 | 995 ± 0.037 | 988 ± 0.038 | 1.003 ± 0.038 | 978 ± 0.038 | 1.009 ± 0.038 | 1.004 ± 0.038 | 1.017 ± 0.038 | 1.012 ± 0.038 |
| | 95 | 840 ± 0.037 | 919 ± 0.039 | 924 ± 0.039 | 918 ± 0.037 | 913 ± 0.037 | 925 ± 0.039 | 938 ± 0.039 | 913 ± 0.039 | 901 ± 0.038 | 886 ± 0.036 | 926 ± 0.037 | 925 ± 0.039 | 925 ± 0.039 | 944 ± 0.038 | 1.001 ± 0.039 | 994 ± 0.038 | 1.055 ± 0.040 | 1.034 ± 0.040 |
| | 90 | 683 ± 0.036 | 835 ± 0.036 | 874 ± 0.040 | 900 ± 0.038 | 871 ± 0.035 | 878 ± 0.039 | 854 ± 0.040 | 832 ± 0.038 | 874 ± 0.037 | 871 ± 0.037 | 846 ± 0.038 | 853 ± 0.037 | 876 ± 0.038 | 870 ± 0.037 | 1.006 ± 0.040 | 1.009 ± 0.041 | 1.108 ± 0.041 | 1.068 ± 0.042 |
| | 85 | 496 ± 0.034 | 763 ± 0.036 | 829 ± 0.039 | 864 ± 0.039 | 894 ± 0.033 | 820 ± | | | | | | | | | | | | |