

Provably Improving Expert Predictions with Prediction Sets

Eleni Straitouri¹, Lequn Wang², Nastaran Okati¹, and Manuel Gomez Rodriguez¹

¹Max Planck Institute for Software Systems, {estraitouri, nastaran, manuelgr}@mpi-sws.org

²Cornell University, lw633@cornell.edu

Abstract

Automated decision support systems promise to help human experts solve tasks more efficiently and accurately. However, existing systems typically require experts to understand when to cede agency to the system or when to exercise their own agency. Moreover, if the experts develop a misplaced trust in the system, their performance may worsen. In this work, we lift the above requirement and develop automated decision support systems that, by design, do not require experts to understand when to trust them to provably improve their performance. To this end, we focus on multiclass classification tasks and consider an automated decision support system that, for each data sample, uses a classifier to recommend a subset of labels to a human expert. We first show that, by looking at the design of such a system from the perspective of conformal prediction, we can ensure that the probability that the recommended subset of labels contains the true label matches almost exactly a target probability value. Then, we develop an efficient and near-optimal search method to find the target probability value under which the expert benefits the most from using our system. Experiments on synthetic and real data demonstrate that our system can help the experts make more accurate predictions and is robust to the accuracy of the classifier it relies on.

1 Introduction

In recent years, there has been an increasing interest in developing automated decision support systems to help human experts solve tasks in a wide range of critical domains, from medicine [1] and drug discovery [2] to criminal justice processes [3], to name a few. Among them, one of the main focus has been classification tasks, where a decision support system uses a classifier to make label predictions and the experts decide when to follow or not follow the predictions made by the classifier [4, 5, 6].

However, these systems typically require the human experts to understand when to trust a prediction made by the classifier. Otherwise, the experts may be better off solving the classification tasks on their own [7]. This follows from the fact that, in general, the accuracy of a classifier differs across data samples [8]. Unfortunately, it is not yet clear how to make sure that the experts do not develop a misplaced trust that decreases their performance [9, 10, 11]. In this work, our goal is to develop decision support systems that, by design, do not require experts to understand when to trust the system and provably improve their performance.

Our contributions. We consider multiclass classification tasks and decision support systems that, for each data sample, use a classifier to recommend a subset of labels to a human expert. We view this type of decision support systems as more natural since, given a set of alternatives, human experts tend to narrow down their options to a subset of them before making their final decision [12, 13, 14]. In a way, our support systems help the experts by automatically narrowing down their options for them. In this context, a recent empirical study has also concluded that, in terms of the overall accuracy, it may be more beneficial to recommend a subset of options than a single option [15].

Here, one could still argue that the expert needs to understand when to trust the system—when to predict one of the labels in the subset recommended by the system—as noted by Levi et al. [15]. This is largely due to the fact that the accuracy of the classifier differs across different data samples. To circumvent this, we

use the theory of conformal prediction [16, 17] to construct *trustworthy* subsets where the probability that the true label belongs to a recommended subset always matches almost exactly a given target probability value, without making any distributional assumptions about the data or the classifier. In addition, given an estimator of the expert’s success probability for any of the recommended subsets, we develop an efficient and near-optimal search method to find the target probability value under which the expert is guaranteed to achieve the greatest accuracy. In this context, we also propose a practical method to obtain such an estimator using the confusion matrix of the expert predictions in the original classification task as well as a given discrete choice model.

Finally, we experiment with synthetic and real data comprising of 511,400 expert predictions over 10,000 natural images. The results demonstrate that our decision support system is robust to the performance of the classifier it relies on—the competitive advantage it provides improves with the accuracy of the classifier, and the human experts do not decrease their performance by using the system even if the classifier is very inaccurate. Additionally, the results also show that, even if the classifiers that our system relies on have high accuracy, an expert using our system may achieve significantly higher accuracy than the classifiers on their own—in our experiments with real data, the relative reduction in misclassification probability is over 61%. Finally, by using our system, our results suggest that the (average) expert would reduce their misclassification probability by $\sim 70\%$.¹

Further related work. Our work builds upon further related work on distribution-free uncertainty quantification and learning under algorithmic triage.

There exist three fundamental notions of distribution-free uncertainty quantification in the literature: calibration, confidence intervals, and prediction sets [17, 18, 19, 16]. Our work is most closely related to the rapidly increasing literature on prediction sets [20, 21, 22, 23], however, to the best of our knowledge, prediction sets have not been optimized to serve automated decision support systems such as ours. In this context, we would like to acknowledge that Babbar et al. [24] have also proposed using prediction sets in decision support systems, however, this work is contemporary to ours and has been carried out independently. Moreover, in contrast to our work, for each data sample, they allow the expert to predict label values outside the recommended subset and do not optimize the probability that the true label belongs to the subset (to provably improve expert predictions).

Learning under algorithmic triage seeks the development of machine learning models that operate under different automation levels—models that take decisions for a given fraction of instances and leave the remaining ones to human experts [8, 25, 26, 27, 28]. This line of work has predominantly focused on supervised learning settings with a few very recent notable exceptions [29, 30]. However, in this line of work, each sample is either predicted by the model or by the human expert. In contrast, in our work, the model helps the human predict each sample.

2 Problem Formulation

We consider a multiclass classification task where a human expert observes a feature vector² $x \in \mathcal{X}$, with $x \sim P(X)$, and needs to predict a label $y \in \mathcal{Y} = \{1, \dots, n\}$, with $y \sim P(Y|X)$. Then, our goal is to design an automated decision support system $\mathcal{C} : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ that, given a feature vector $x \in \mathcal{X}$, helps the expert by automatically narrowing down the set of potential labels to a subset of them $\mathcal{C}(x) \subseteq \mathcal{Y}$ using a trained classifier $\hat{f}(x) \in [0, 1]^n$ that outputs scores for each class (*e.g.*, softmax scores)³. The higher the score $\hat{f}_y(x)$, the more the classifier believes the true label $Y = y$. Here, we assume that, for each $x \sim P(X)$, the expert predicts a label \hat{Y} among those in the subset $\mathcal{C}(x)$ according to an unknown policy $\pi(x, \mathcal{C}(x))$. More formally, $\hat{Y} \sim \pi(x, \mathcal{C}(x))$, where $\pi : \mathcal{X} \times 2^{\mathcal{Y}} \rightarrow \Delta(\mathcal{Y})$ and $\Delta(\mathcal{Y})$ denotes the probability simplex over the set of labels \mathcal{Y} , and $\pi_y(x, \mathcal{C}(x)) = 0$ if $y \notin \mathcal{C}(x)$. Refer to Figure 1 for an illustration of the automated decision support system we consider.

¹To facilitate research in this area, we will release an open-source implementation of our system with the final version of the paper.

²We denote random variables with capital letters and realizations of random variables with lower case letters.

³The assumption that $\hat{f}(x) \in [0, 1]^n$ is without loss of generality.

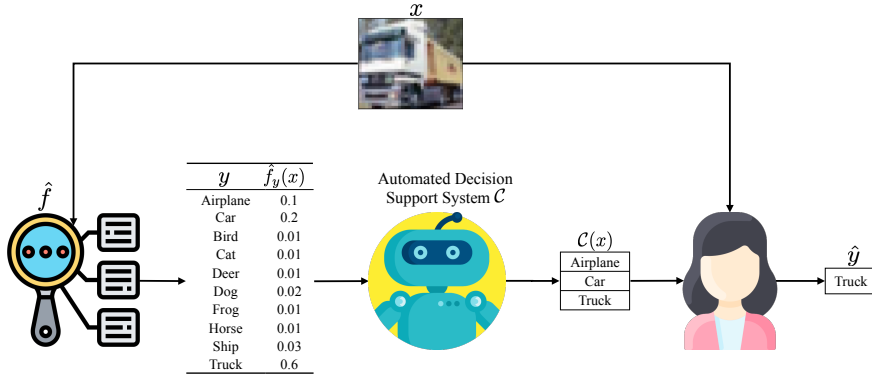


Figure 1: Our automated decision support system \mathcal{C} . Given a sample with a feature vector x , our system \mathcal{C} narrows down the set of potential labels $y \in \mathcal{Y}$ to a subset of them $\mathcal{C}(x)$ using the scores $\hat{f}_y(x)$ provided by a classifier \hat{f} for each class y . The human expert receives the recommended subset $\mathcal{C}(x)$, together with the sample, and predicts a label \hat{y} from $\mathcal{C}(x)$ according to a policy $\pi(x, \mathcal{C}(x))$.

In our work, we would like that, by design, the expert can only benefit from using the automated decision support system \mathcal{C} , *i.e.*,

$$\mathbb{P}[\hat{Y} = Y ; \mathcal{C}] \geq \mathbb{P}[\hat{Y} = Y ; \mathcal{Y}], \quad (1)$$

where $\mathbb{P}[\hat{Y} = Y ; \mathcal{C}]$ denotes the expert’s success probability if, for each $x \sim P(X)$, the human expert predicts a label \hat{Y} among those in the subset $\mathcal{C}(x)$.

However, not all automated decision support systems fulfilling the above requirement will be equally useful—some will help experts increase their success probability more than others. For example, a system that always recommends $\mathcal{C}(x) = \mathcal{Y}$ for all $x \in \mathcal{X}$ satisfies Eq. 1. However, it is useless to the experts. Therefore, among those systems satisfying Eq. 1, we would ideally like to find the system \mathcal{C}^* that helps the experts achieve the highest success probability⁴, *i.e.*,

$$\mathcal{C}^* = \operatorname{argmax}_{\mathcal{C}} \mathbb{P}[\hat{Y} = Y ; \mathcal{C}]. \quad (2)$$

To address the design of such a system, we will look at the problem from the perspective of conformal prediction [17, 16].

3 Subset Selection using Conformal Prediction

In general, if the trained classifier \hat{f} we use to build $\mathcal{C}(X)$ is not perfect, the true label Y may or may not be included in $\mathcal{C}(X)$. In what follows, we will construct the subsets $\mathcal{C}(X)$ using the theory of conformal prediction. This will allow our system to be robust to the accuracy of the classifier \hat{f} it uses—the probability $P[Y \in \mathcal{C}(X)]$ that the true label Y belongs to the subset $\mathcal{C}(X) = \mathcal{C}_\alpha(X)$ will match almost exactly a given target probability value $1 - \alpha$, without making any distributional assumptions about the data distribution $P(X)P(Y | X)$ nor the classifier \hat{f} .

Let $\mathcal{D}_{\text{cal}} = \{(x_i, y_i)\}_{i=1}^m$ be a calibration set, where $(x_i, y_i) \sim P(X)P(Y | X)$, $s(x_i, y_i) = 1 - \hat{f}_{y_i}(x_i)$ be the *conformal score*⁵ (*i.e.*, if the classifier is catastrophically wrong, the conformal score will be close to one), and \hat{q}_α be the $\frac{\lceil (m+1)(1-\alpha) \rceil}{m}$ empirical quantile of the conformal scores $s(x_1, y_1), \dots, s(x_m, y_m)$. Then, if we

⁴Note that maximizing the expert’s success probability $\mathbb{P}[\hat{Y} = Y ; \mathcal{C}]$ is equivalent to minimizing the expected 0-1 loss $\mathbb{E}[\mathbb{I}(\hat{Y} \neq Y) ; \mathcal{C}]$. Considering other types of losses is left as an interesting venue for future work.

⁵In general, the conformal score $s(x, y)$ can be any function of x and y measuring the *similarity* between samples.

use the quantile \hat{q}_α to construct the subsets $C_\alpha(X)$ for new data samples as follows:

$$C_\alpha(X) = \{y \mid s(X, y) \leq \hat{q}\}, \quad (3)$$

then it holds that the probability that the true label Y belongs to the subset $C_\alpha(X)$ is almost exactly $1 - \alpha$. More formally, we have the following well-established conformal calibration coverage guarantee (refer to Appendix A in Angelopoulos and Bates [16] for a proof):

Theorem 1 *For an automated decision support system C_α that constructs the subsets $C_\alpha(X)$ using Eq. 3, it holds that*

$$1 - \alpha \leq \mathbb{P}[Y \in C_\alpha(X)] \leq 1 - \alpha + \frac{1}{m+1},$$

where the probability is over the randomness in the sample it helps predicting and the calibration set used to compute the empirical quantile \hat{q}_α .

However, we would like to emphasize that, given a fixed calibration set \mathcal{D}_{cal} , there may be some α values that will lead to larger gains in terms of success probability $\mathbb{P}[\hat{Y} = Y; C_\alpha]$ than others. In this context, note that the success probability is only over the randomness in the samples the system helps predicting—the calibration set \mathcal{D}_{cal} is fixed. This is in contrast with the coverage guarantee given by Theorem 1, where the probability is both over the randomness in the samples the system helps predicting and the calibration set. In what follows, our goal is to find the optimal value α^* that maximizes the expert’s success probability given a fixed calibration set \mathcal{D}_{cal} .

4 Optimizing Across Conformal Predictors

We start by realizing that, given a fixed calibration set $\mathcal{D}_{\text{cal}} = \{(x_i, y_i)\}_{i=1}^m$, there only exist m different conformal predictors. This is because the empirical quantile \hat{q}_α , which the subsets $C_\alpha(x_i)$ depend on, can only take m different values. As a result, to find the optimal conformal predictor that maximizes the expert’s success probability, we need to solve the following maximization problem:

$$\alpha^* = \operatorname{argmax}_{\alpha \in \mathcal{A}} \mathbb{P}[\hat{Y} = Y; C_\alpha], \quad (4)$$

where $\mathcal{A} = \{\alpha_i\}_{i \in [m]}$, with $\alpha_i = 1 - i/(m+1)$, and the probability is only over the randomness in the samples the system helps predicting.

However, to find a near optimal solution $\hat{\alpha}$ to the above problem, we need to estimate the expert’s success probability $\mathbb{P}[\hat{Y} = Y; C_\alpha]$. Assume for now that, for each $\alpha \in \mathcal{A}$, we have access to an estimator $\hat{\mu}_\alpha$ of the expert’s success probability such that, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, it holds that $|\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; C_\alpha]| \leq \epsilon_{\alpha, \delta}$. Then, we can use the following proposition to find a near optimal solution $\hat{\alpha}$ to Eq. 4 with high probability:

Proposition 1 *For any $\delta \in (0, 1)$, consider an automated decision support system $C_{\hat{\alpha}}$ with*

$$\hat{\alpha} = \operatorname{argmax}_{\alpha \in \mathcal{A}} \hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m}. \quad (5)$$

With probability at least $1 - \delta$, it holds that $\mathbb{P}[\hat{Y} = Y; C_{\hat{\alpha}}] \geq \mathbb{P}[\hat{Y} = Y; C_\alpha] - 2\epsilon_{\alpha, \delta/m} \forall \alpha \in \mathcal{A}$ simultaneously.

More specifically, the above result directly implies that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, it holds that:

$$\mathbb{P}[\hat{Y} = Y; C_{\alpha^*}] - \mathbb{P}[\hat{Y} = Y; C_{\hat{\alpha}}] \leq 2\epsilon_{\alpha^*, \delta/m}. \quad (6)$$

In what follows, we propose a practical method to estimate the expert’s success probability $\mathbb{P}[\hat{Y} = y; C_\alpha]$ that builds upon the multinomial logit model (MNL), one of the most popular models in the vast literature

Algorithm 1 Finding a near-optimal $\hat{\alpha}$

```
1: Input:  $\hat{f}$ ,  $\mathcal{D}_{\text{est}}$ ,  $\mathcal{D}_{\text{cal}}$ ,  $\delta$ 
2: Initialize:  $\mathcal{A} = \{\}$ ,  $\hat{\alpha} \leftarrow 0$ ,  $t \leftarrow 0$ 
3: for  $i \in 1, \dots, m$  do
4:    $\alpha \leftarrow 1 - \frac{i}{m+1}$ 
5:    $\mathcal{A} \leftarrow \mathcal{A} \cup \{\alpha\}$ 
6: end for
7: for  $\alpha \in \mathcal{A}$  do
8:    $\hat{\mu}_\alpha, \epsilon_{\alpha, \delta/m} \leftarrow \text{ESTIMATE}(\alpha, \delta, \mathcal{D}_{\text{est}})$  {It uses Eqs. 8 and 9}
9:   if  $t \leq \hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m}$  then
10:     $t \leftarrow \hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m}$ 
11:     $\hat{\alpha} \leftarrow \alpha$ 
12:   end if
13: end for
14: return  $\hat{\alpha}$ 
```

on discrete choice models [31]. More specifically, we assume that we have access to (an estimation of) the confusion matrix \mathbf{C} for the expert predictions in the (original) multiclass classification task, similarly as in Kerrigan et al. [32], *i.e.*,

$$\mathbf{C} = [C_{yy'}]_{y, y' \in \mathcal{Y}}, \text{ where } C_{yy'} = \mathbb{P}[\hat{Y} = y' \mid \mathcal{Y}, Y = y].$$

Moreover, given a sample (x, y) , we also assume that the expert's conditional success probability for the subset $\mathcal{C}_\alpha(x)$ follows a multinomial logit model (MNL) [31], *i.e.*,

$$\mathbb{P}[\hat{Y} = y; \mathcal{C}_\alpha \mid y \in \mathcal{C}_\alpha(x)] = \frac{e^{u_{yy}}}{\sum_{y' \in \mathcal{C}_\alpha(x)} e^{u_{yy'}}}, \quad (7)$$

where $u_{yy'} = \log C_{yy'}$. Then, building on the above expression, we compute a Monte-Carlo estimator $\hat{\mu}_\alpha$ of the conditional success probability using an estimation set $\mathcal{D}_{\text{est}} = \{(x_i, y_i)\}_{i \in [m]}$ ⁶, *i.e.*,

$$\hat{\mu}_\alpha = \frac{1}{m} \sum_{i \in [m] \mid y_i \in \mathcal{C}_\alpha(x_i)} \mathbb{P}[\hat{Y} = y_i \mid \mathcal{C}_\alpha(x_i), y_i \in \mathcal{C}_\alpha(x_i)]. \quad (8)$$

Here, for each $\alpha \in \mathcal{A}$, using Hoeffding's inequality⁷, we can conclude that, with probability at least $1 - \delta$, it holds that (refer to Appendix B):

$$|\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]| \leq \sqrt{\frac{\log \frac{1}{\delta}}{2m}} := \epsilon_{\alpha, \delta}. \quad (9)$$

As a consequence, as $m \rightarrow \infty$, $\epsilon_{\alpha, \delta}$ converges to zero.

Algorithm 1 summarizes the overall search method, which first builds \mathcal{A} and then finds the near optimal $\hat{\alpha}$ in \mathcal{A} . To build \mathcal{A} , it needs $\mathcal{O}(m)$ steps. To find the near-optimal $\hat{\alpha}$, for each value $\alpha \in \mathcal{A}$ and each sample $(x_i, y_i) \in \mathcal{D}_{\text{est}}$, it needs to compute a prediction set $\mathcal{C}_\alpha(x)$. This is achieved by sorting the conformal scores and reusing computations across α values, which takes $\mathcal{O}(m \log m + mn \log n)$ steps. Therefore, the overall time complexity is $\mathcal{O}(m \log m + mn \log n)$.

5 Beyond Standard Conformal Prediction

Until now, we have used standard conformal prediction [16] to construct the recommended subsets $\mathcal{C}(X)$ —we have constructed the subsets $\mathcal{C}(X)$ by comparing the conformal scores $s(X, y)$ to a single threshold \hat{q} , as

⁶The number of samples in \mathcal{D}_{cal} and \mathcal{D}_{est} can differ. However, for simplicity, we assume both sets contain m samples.

⁷We use Hoeffding's inequality for ease of exposition, however, one could use tighter concentration inequalities.

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.41	0.57	0.74	0.91
0.5	0.55	0.68	0.81	0.93
0.7	0.72	0.79	0.88	0.96
0.9	0.90	0.92	0.95	0.98

Table 1: Average success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}_1, \hat{\alpha}_2}]$ during test for four different experts using our system, each with a different success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, on four prediction tasks where the classifier achieves a different success probability $\mathbb{P}[Y' = Y]$. Each column corresponds to a different prediction task and each row corresponds to a different expert. In each task, the number of label values $n = 10$ and the size of the calibration and estimation sets is $m = 1,200$. Each cell shows only the average since the standard errors are all below 10^{-2} .

shown in Eq. 3. Here, we will show that we can sometimes improve the performance of our system by constructing $\mathcal{C}(X)$ using two thresholds \hat{q}_{α_1} and \hat{q}_{α_2} . By doing so, the recommended subsets will include label values whose corresponding conformal scores are neither unreasonably large, as in standard conformal prediction, nor unreasonably low, in comparison with the conformal scores of the samples in the calibration set \mathcal{D}_{cal} . This will be particularly useful whenever the classifier our system relies on has low accuracy.

More specifically, given a calibration set $\mathcal{D}_{\text{cal}} = \{(x_i, s_i)\}_{i=1}^m$, let $\alpha_1, \alpha_2 \in [0, 1]$, with $\alpha_1 < \alpha_2$, and \hat{q}_{α_1} and \hat{q}_{α_2} be the $\frac{\lceil (m+1)(1-\alpha_1) \rceil}{m}$ and $\frac{\lceil (m+1)(1-\alpha_2) \rceil}{m}$ empirical quantiles of the conformal scores $s(x_1, y_1), \dots, s(x_m, y_m)$. Then, if we use the quantiles \hat{q}_{α_1} and \hat{q}_{α_2} to construct the subsets $\mathcal{C}_{\alpha_1, \alpha_2}(X)$ for new data samples as follows:

$$\mathcal{C}_{\alpha_1, \alpha_2}(X) = \{y \mid \hat{q}_{\alpha_2} < s(X, y) \leq \hat{q}_{\alpha_1}\}, \quad (10)$$

it holds that the probability that the true label Y belongs to the subset $\mathcal{C}_{\alpha_1, \alpha_2}(X)$ is almost exactly $\alpha_2 - \alpha_1$. More formally, we have the following conformal calibration coverage guarantee, which is the counterpart of Theorem 1:

Theorem 2 Consider an automated decision support system $\mathcal{C}_{\alpha_1, \alpha_2}$ that constructs $\mathcal{C}_{\alpha_1, \alpha_2}(X)$ using Eq. 10. Then, it holds that:

$$\alpha_2 - \alpha_1 - \frac{1}{m+1} \leq \mathbb{P}[Y \in \mathcal{C}_{\alpha_1, \alpha_2}(X)] \leq \alpha_2 - \alpha_1 + \frac{1}{m+1}.$$

Moreover, given an estimator of the expert’s success probability $\hat{\mu}_{\alpha_1, \alpha_2}$ such that for each $\alpha_1 < \alpha_2$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$, it holds that $|\hat{\mu}_{\alpha_1, \alpha_2} - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\alpha_1, \alpha_2}]| \leq \epsilon_{\alpha_1, \alpha_2, \delta}$, we can proceed similarly as in standard conformal prediction to find the near optimal $\hat{\alpha}_1, \hat{\alpha}_2 \in \mathcal{A}$ that maximizes the expert’s success probability with high probability, by using $\hat{\mu}_{\alpha_1, \alpha_2}$ and $\epsilon_{\alpha_1, \alpha_2, 2\delta/(m(m-1))}$. Here, it is worth pointing out that, in contrast with the case of standard conformal prediction, the time complexity of finding the near optimal $\hat{\alpha}_1$ and $\hat{\alpha}_2$ is $\mathcal{O}(m \log m + mn \log n + mn^2)$. Moreover, we can still rely on the practical method to estimate the expert’s conditional success probability introduced in Section 4.

Remarks. Conformal prediction is one of many possible ways to construct set-valued predictors [33], *i.e.*, predictors that, for each sample $x \in \mathcal{X}$, output a set of label candidates $\mathcal{C}(x)$. In our work, we favor conformal predictors over alternatives because they provably output *trustworthy* sets $\mathcal{C}_\alpha(x)$ without making any assumption about the data distribution nor the classifier they rely upon. However, we would like to emphasize that our efficient search method (Algorithm 1) is rather generic and, together with an estimator of the expert’s success probability, may be used to find a near-optimal set-valued predictor within a discrete set of set-valued predictors that maximizes the expert’s success probability. We hope our work will encourage others to develop set-valued predictors specifically designed to serve decision support systems.

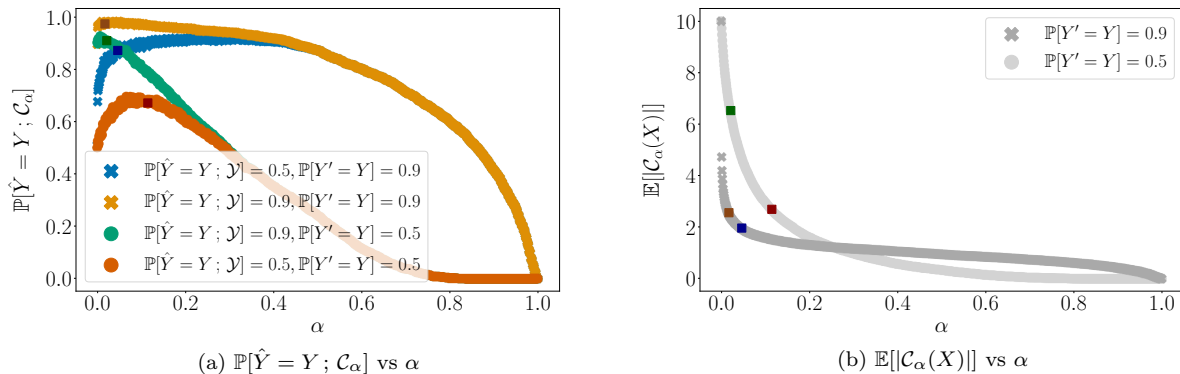


Figure 2: Expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]$ and average size of the recommended sets $\mathbb{E}[|\mathcal{C}_\alpha(X)|]$ during test for each $\alpha \in \mathcal{A}$ on four synthetic prediction tasks where the expert and the classifier achieve different success probabilities $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$ and $\mathbb{P}[Y' = Y]$, respectively. Here, note that $\mathbb{E}[|\mathcal{C}_\alpha(X)|]$ only depends on the classifier’s success probability $\mathbb{P}[Y' = Y]$ and thus we only need two lines for the four prediction tasks. In all tasks, the number of label values $n = 10$ and the size of the calibration and estimation sets is $m = 1,200$. Each marker corresponds to a different α value and the darker points correspond to $\hat{\alpha}$. The coloring of the darker points for each prediction task is the same in both panels.

6 Experiments on Synthetic Data

In this section, we quantify the influence that the accuracy of the expert and the classifier, the size of the calibration and estimation sets, as well as the number of label values have on the performance of our system using synthetic data⁸.

Experimental setup. We create a variety of synthetic prediction tasks, each with 20 features per sample and a varying number of label values n and difficulty. Refer to Appendix D for more details about the prediction tasks. For each prediction task, we generate 10,000 samples, pick 20% of these samples at random as test set and split also at random the remaining 80% into three disjoint subsets for training, calibration and estimation, whose sizes we vary across experiments. In each experiment, we specify the number of samples in the calibration and estimation sets—the remaining samples are used for training.

For each prediction task, we train a logistic regression model $P_\theta(Y' | X)$, which depending on the difficulty of the prediction task, achieves different success probability values $\mathbb{P}[Y' = Y]$. Moreover, we sample the expert’s predictions \hat{Y} from the multinomial logit model defined by Eq. 7, with $C_{yy} = \frac{\pi}{n} \pm \gamma \epsilon_c$ and $C_{yy'} = \frac{1 - C_{yy}}{n} \pm \beta$, where π is a parameter that controls the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, $\epsilon_c \sim U(0, \min(1 - \frac{\pi}{n}, \frac{\pi}{n}))$, $\beta \sim N(0, ((1 - C_{yy})/(6n))^2)$ for all $y \neq y'$, and γ is a normalization term. Finally, we repeat each experiment five times and, each time, we sample different train, estimation, calibration and test sets following the procedure described above.

Experts always benefit from our system even if the classifier has low accuracy. We estimate the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}_1, \hat{\alpha}_2}]$ during test for four different experts, each with a different success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, on four prediction tasks where the classifier achieves a different success probability $\mathbb{P}[Y' = Y]$. Table 1 summarizes the results, where each column corresponds to a different prediction task and each row corresponds to a different expert. We find that, using our system, the expert solves the prediction task significantly more accurately than the expert or the classifier on their own. Moreover, it is rather remarkable that, even if the classifier has low accuracy, the expert always benefits from using our system—in other words, our system is robust to the performance of the classifier it relies on. We found

⁸All algorithms ran on a Debian machine equipped with Intel Xeon E5-2667 v4 @ 3.2 GHz, 32GB memory and two M40 Nvidia Tesla GPU cards. See Appendix D for further details.

CLASSIFIERS	$\mathbb{P}[Y' = Y]$	$\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$
RESNET-110	0.928	0.981
PRERESNET-110	0.944	0.983
DENSENET	0.964	0.986

Table 2: Average success probabilities $\mathbb{P}[Y' = Y]$ and $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ achieved by three popular deep neural network classifiers and by an expert using our system with each of the classifiers on the CIFAR-10H dataset. The size of the calibration and estimation set is $m = 1,500$ and the expert’s average success probability at solving the (original) multiclass task is $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}] = 0.947$. Each cell shows only the average over 10 random data splits since the standard errors are all below 10^{-2} .

qualitatively similar results for prediction tasks with other values of n and m , which we report in Appendix E. Since, except for prediction tasks where the classifier has the lowest accuracy (*e.g.*, $\mathbb{P}[Y' = Y] = 0.3$), we found that typically $\hat{\alpha}_2 = 1$, in what follows, we only experiment with systems $\mathcal{C}_{\hat{\alpha}} = \mathcal{C}_{\hat{\alpha},1}$ that construct $\mathcal{C}_{\hat{\alpha}}(X)$ using Eq. 3.

The performance of our system under $\hat{\alpha}$ found by Algorithm 1 and under α^* is very similar. Given three prediction tasks where the expert and the classifier achieve different success probabilities $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$ and $\mathbb{P}[Y' = Y]$, we compare the performance of our system under the near optimal $\hat{\alpha}$ found by Algorithm 1 and under all other possible $\alpha \in \mathcal{A}$ values. Figure 2 summarizes the results, which suggest that: (i) the performance under $\hat{\alpha}$ is very close to that under α^* , as suggested by Proposition 1; and, (ii) as long as $\alpha \leq \alpha^*$, the performance of our system increases monotonically with respect to α , however, once $\alpha > \alpha^*$, the performance deteriorates as we increase α . (iii) the higher the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, the smaller the optimal α^* and thus the greater the average size of the subsets $\mathcal{C}_{\alpha^*}(X)$. We found qualitatively similar results using other expert-classifier pairs with different success probabilities.

Our system needs a relatively small amount of calibration and estimation data. We vary the amount of calibration and estimation data m we feed into Algorithm 1 and, each time, estimate the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$. Across prediction tasks, we consistently find that our system needs a relatively small amount of calibration and estimation data to perform well. For example, for all prediction tasks with $n = 10$ label values and varying level of difficulty, if we increase the amount of calibration and estimation data from $m = 160$ to $m = 1,200$, the relative gain in success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ with respect to $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, averaged across experts with $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}] \in \{0.3, 0.5, 0.7, 0.9\}$, goes from $47.29 \pm 6.37\%$ to $48.81 \pm 6.65\%$.

The greater the number of label values, the more an expert benefits from using our system. We consider prediction tasks with a varying number of label values, from $n = 10$ to $n = 100$, and estimate the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ for each task. Our results suggest that the relative gain in success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ with respect to $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, averaged across experts with $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}] \in \{0.3, 0.5, 0.7, 0.9\}$, increases with the number of label values. For example, for $m = 400$, it goes from $48.04 \pm 6.69\%$ for $n = 10$ to $71.32 \pm 9.30\%$ for $n = 100$. For other m values, we found a similar trend.

7 Experiments on Real Data

In this section, we evaluate the performance of our system using a dataset with real expert predictions over natural images and several popular and highly accurate deep neural network classifiers. Here, we focus on systems $\mathcal{C}_{\alpha} = \mathcal{C}_{\alpha,1}$ that construct $\mathcal{C}_{\alpha}(X)$ using Eq. 3 because, whenever the classifiers are highly accurate, systems $\mathcal{C}_{\alpha_1, \alpha_2}$ with $\alpha_2 \neq 1$ do not offer a competitive advantage.

Data description and experimental setup. We experiment with the dataset CIFAR-10H [34], which contains 10,000 natural images taken from the test set of the standard CIFAR-10 [35]. Each of these images

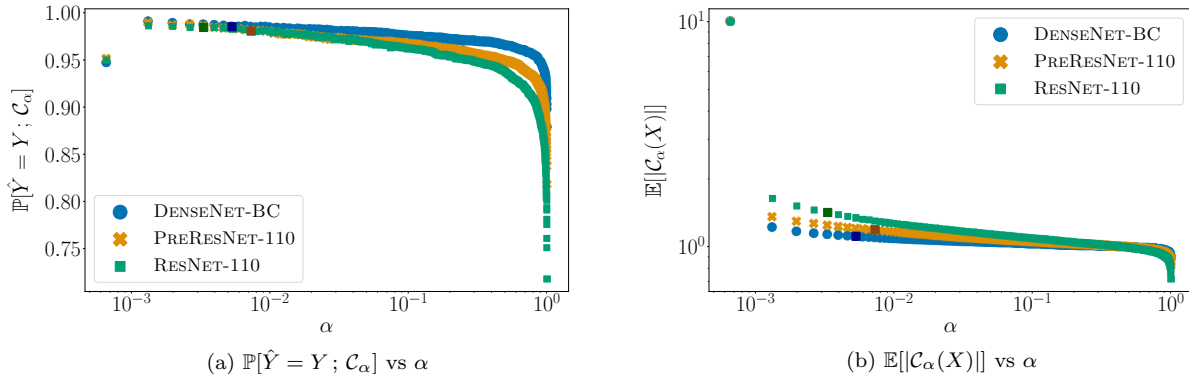


Figure 3: Expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]$ and average size of the recommended sets $\mathbb{E}[|\mathcal{C}_\alpha(X)|]$ during test for each $\alpha \in \mathcal{A}$ for three popular deep neural network classifiers on the CIFAR-10H dataset. The size of the calibration and estimation sets is $m = 1,500$. Each marker corresponds to a different α value and the darker points correspond to $\hat{\alpha}$ for each task.

belongs to one of $n = 10$ classes and contains approximately 50 expert predictions \hat{Y}^9 . Here, we randomly split the dataset into three disjoint subsets for calibration, estimation and test, whose sizes we vary across experiments. In each experiment, we specify the number of samples in the calibration and estimation sets—the remaining samples are used for testing.

Rather than training a classifier, we use three popular and highly accurate deep neural network classifiers trained on CIFAR-10, namely ResNet-110 [36], PreResNet-110 [37] and DenseNet [38]. Moreover, we use the human predictions \hat{Y} to estimate the confusion matrix \mathbf{C} for the expert predictions in the (original) multiclass classification task [32] and then sample the expert’s prediction \hat{Y} from the multinomial logit model defined by Eq. 7 to both estimate the expert’s conditional success probabilities in Eq. 8 in Algorithm 1 and estimate the expert’s success probability during test.

Results. We start by estimating the success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ during test for an expert using our system with each of the above mentioned classifiers. Table 2 summarizes the results, where we compare the success probability achieved by our system and by the corresponding classifier on its own and we report the expert’s success probability at solving the (original) multiclass task in the caption. We find it very encouraging that, despite the classifiers are highly accurate, our results suggest that an expert using our system can solve the prediction task significantly more accurately than the classifiers. More specifically, the relative reduction in misclassification probability goes from 61.1% (DenseNet) to a striking 73.6% (ResNet-110). Moreover, by using our system, our results suggest that the (average) expert would reduce their misclassification probability by $\sim 70\%$.

Next, for each choice of classifier, we compare the performance of our system under the near optimal $\hat{\alpha}$ found by Algorithm 1 and under all other possible α values, including the optimal α^* . Figure 3 summarizes the results, which suggest that, similarly as in the experiments on synthetic data, the performance of our system under $\hat{\alpha}$ and α^* is very similar. However, since the classifiers are all highly accurate, the average size of the recommended subsets under $\hat{\alpha}$ and α^* is quite close to one even though $\hat{\alpha}$ is much smaller than in the experiments in synthetic data.

Finally, we also investigate to what extent the amount of calibration and estimation data m we feed into Algorithm 1 influences the expert’s success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}]$ under our system. In contrast with the experiments on synthetic data, we do find that our system needs larger amounts of calibration and estimation

⁹The dataset CIFAR-10H is among the only publicly available datasets that we found containing multiple expert predictions per sample, necessary to estimate the confusion matrix \mathbf{C} , a relatively large number of samples, and more than two classes. The dataset is released under Creative Commons BY-NC-SA 4.0 license. However, since our methodology is rather general, our system may help improving expert predictions in other applications.

data to realize its full potential. For example, while the relative gain in success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]$ with respect to $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$ is just $1.51 \pm 0.3\%$ under $m = 200$, it raises to $2.99 \pm 0.15\%$ under $m = 1,500$. However, this is not surprising since the classifiers are very accurate and thus the optimal α^* is very small. Therefore, we need larger calibration sets to have enough granularity to distinguish among small α values.

8 Conclusions

We have initiated the development of automated decision support systems that, by design, do not require human experts to understand when to trust them to provably improve their performance. In particular, we have focused on multiclass classification tasks and designed a system that, for each data sample, recommends a subset of labels to the experts using a classifier. Moreover, we have shown that our system can help experts make predictions more accurately and is robust to the performance of the classifier it relies on.

Our work opens up many interesting avenues for future work. For example, we have considered a simple conformal score function from the literature. It would be interesting to develop score functions especially designed to improve the performance of our system. Moreover, to estimate the expert’s success probability, we have introduced a simple procedure that builds upon the multinomial logit model (MNL), a classical model from the discrete choice model literature. However, more research is needed to accurately estimate the expert’s success probability, *e.g.*, using more sophisticated discrete choice models. In this context, it would be also interesting to develop systems that perform online estimation of the expert’s conditional success probability. In addition, it would be important to investigate the ethical impact of our decision support system, including human trust and bias. In this context, we hypothesize that our system may sometimes help reduce human bias, *e.g.*, in cases in which the provided sets exclude a label that a human systematically mistakes as another. However, one can also think of counterexamples where our system may introduce new sources of bias. Finally, it would be important to deploy and evaluate our system on a real-world application with human experts.

Acknowledgements. Gomez-Rodriguez acknowledges support from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 945719).

References

- [1] Wei Jiao, Gurnit Atwal, Paz Polak, Rosa Karlic, Edwin Cuppen, Alexandra Danyi, Jeroen de Ridder, Carla van Herpen, Martijn P Lolkema, Neeltje Steeghs, et al. A deep learning system accurately classifies primary and metastatic cancers using passenger mutation patterns. *Nature communications*, 11(1):1–12, 2020.
- [2] Ruishan Liu, Shemra Rizzo, Samuel Whipple, Navdeep Pal, Arturo Lopez Pineda, Michael Lu, Brandon Arnieri, Ying Lu, William Capra, Ryan Copping, et al. Evaluating eligibility criteria of oncology trials using real-world data and ai. *Nature*, 592(7855):629–633, 2021.
- [3] Nina Grgić-Hlača, Christoph Engel, and Krishna P Gummadi. Human decision making with machine assistance: An experiment on bailing and jailing. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–25, 2019.
- [4] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S. Lasecki, Daniel S. Weld, and Eric Horvitz. Beyond accuracy: The role of mental models in human-ai team performance. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 7(1):2–11, Oct. 2019.
- [5] Brian Lubars and Chenhao Tan. Ask not what ai can do, but what ai should do: Towards a framework of task delegability. *Advances in Neural Information Processing Systems*, 32:57–67, 2019.
- [6] Sebastian Bordt and Ulrike von Luxburg. When humans and machines make joint decisions: A non-symmetric bandit model. *arXiv preprint arXiv:2007.04800*, 2020.

- [7] Harini Suresh, Natalie Lao, and Ilaria Liccardi. Misplaced trust: Measuring the interference of machine learning in human decision-making. In *12th ACM Conference on Web Science*, pages 315–324, 2020.
- [8] Maithra Raghu, Katy Blumer, Greg Corrado, Jon Kleinberg, Ziad Obermeyer, and Sendhil Mullainathan. The algorithmic automation problem: Prediction, triage, and human effort. *arXiv preprint arXiv:1903.12220*, 2019.
- [9] Mahsan Nourani, Joanie T. King, and Eric D. Ragan. The role of domain expertise in user trust and the impact of first impressions with intelligent systems. *ArXiv*, abs/2008.09100, 2020.
- [10] Ming Yin, Jennifer Wortman Vaughan, and Hanna Wallach. Understanding the effect of accuracy on trust in machine learning models. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–12, 2019.
- [11] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. Effect of confidence and explanation on accuracy and trust calibration in ai-assisted decision making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 295–305, 2020.
- [12] Peter Wright and Fredrick Barbour. *Phased decision strategies: Sequels to an initial screening*. Graduate School of Business, Stanford University, 1977.
- [13] Lee Roy Beach. Broadening the definition of decision making: The role of prechoice screening of options. *Psychological Science*, 4(4):215–220, 1993.
- [14] Moshe Ben-Akiva and Bruno Boccara. Discrete choice models with latent choice sets. *International Journal of Research in Marketing*, 12(1):9–24, 1995.
- [15] Ariel Levy, Monica Agrawal, Arvind Satyanarayan, and David Sontag. Assessing the impact of automated suggestions on decision making: Domain experts mediate model errors but take less initiative. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021.
- [16] Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- [17] Vladimir Vovk, Alex Gammerman, and Glenn Shafer. *Algorithmic Learning in a Random World*. Springer-Verlag, Berlin, Heidelberg, 2005.
- [18] Vineeth Balasubramanian, Shen-Shyang Ho, and Vladimir Vovk. *Conformal prediction for reliable machine learning: theory, adaptations and applications*. Newnes, 2014.
- [19] Chirag Gupta, Aleksandr Podkopaev, and Aaditya Ramdas. Distribution-free binary classification: prediction sets, confidence intervals and calibration. In *Advances in Neural Information Processing Systems*, 2020.
- [20] Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. *Advances in Neural Information Processing Systems*, 32:3543–3553, 2019.
- [21] Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33:3581–3591, 2020.
- [22] Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I Jordan. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Representations*, 2021.
- [23] Aleksandr Podkopaev and Aaditya Ramdas. Distribution-free uncertainty quantification for classification under label shift. In *Proceedings of the 37th Conference on Uncertainty in Artificial Intelligence*, 2021.

- [24] Varun Babbar, Umang Bhatt, and Adrian Weller. On the utility of prediction sets in human-ai teams. *arXiv preprint arXiv:2205.01411*, 2022.
- [25] Abir De, Paramita Koley, Niloy Ganguly, and Manuel Gomez-Rodriguez. Regression under human assistance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [26] Abir De, Nastaran Okati, Ali Zarezade, and Manuel Gomez-Rodriguez. Classification under human assistance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [27] Hussein Mozannar and David Sontag. Consistent estimators for learning to defer to an expert. In *International Conference on Machine Learning*, pages 7076–7087, 2020.
- [28] Nastaran Okati, Abir De, and Manuel Gomez-Rodriguez. Differentiable learning under triage. In *Advances in Neural Information Processing Systems*, 2021.
- [29] Vahid Balazadeh Meresht, Abir De, Adish Singla, and Manuel Gomez-Rodriguez. Learning to switch between machines and humans. *arXiv preprint arXiv:2002.04258*, 2020.
- [30] Eleni Straitouri, Adish Singla, Vahid Balazadeh Meresht, and Manuel Gomez-Rodriguez. Reinforcement learning under algorithmic triage. *arXiv preprint arXiv:2109.11328*, 2021.
- [31] Florian Heiss. *Discrete choice methods with simulation*. Taylor & Francis, 2016.
- [32] Gavin Kerrigan, Padhraic Smyth, and Mark Steyvers. Combining human predictions with model probabilities via confusion matrices and calibration. *arXiv preprint arXiv:2109.14591*, 2021.
- [33] Evgenii Chzhen, Christophe Denis, Mohamed Hebiri, and Titouan Lorieul. Set-valued classification–overview via a unified framework. *arXiv preprint arXiv:2102.12318*, 2021.
- [34] Joshua C. Peterson, Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. Human uncertainty makes classification more robust. *arXiv preprint arXiv:1908.07086*, 2019.
- [35] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Citeseer*, 2009.
- [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [38] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, 2017.

A Proof of Proposition 1

Given the estimators $\hat{\mu}_\alpha$ of $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]$, we have that, for each $\alpha \in \mathcal{A}$, it holds that

$$\left| \hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \right| \leq \epsilon_{\alpha, \delta/m} \quad (11)$$

with probability at least $1 - \delta/m$. By applying the union bound, we know that the above events hold simultaneously for all $\alpha \in \mathcal{A}$ with probability at least $1 - \delta$. Moreover, by rearranging, the above expression can be rewritten as

$$\hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m} \leq \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \leq \hat{\mu}_\alpha + \epsilon_{\alpha, \delta/m}. \quad (12)$$

Let $\hat{\alpha} = \operatorname{argmax}_{\alpha \in \mathcal{A}} \{(\hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m})\}$. For $\hat{\alpha}$, with probability $1 - \delta$, it holds that for all $\alpha \in \mathcal{A}$,

$$\begin{aligned} \mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}}] &\geq \hat{\mu}_{\hat{\alpha}} - \epsilon_{\hat{\alpha}, \delta/m} \geq \hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m} = \hat{\mu}_\alpha - \epsilon_{\alpha, \delta/m} + 2\epsilon_{\alpha, \delta/m} - 2\epsilon_{\alpha, \delta/m} \\ &= \hat{\mu}_\alpha + \epsilon_{\alpha, \delta/m} - 2\epsilon_{\alpha, \delta/m}, \\ &\geq \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] - 2\epsilon_{\alpha, \delta/m}, \end{aligned}$$

where the last inequality follows from Eq. 12.

B Derivation of error expression for Hoeffding's inequality

From Hoeffding's inequality we have that:

Theorem 3 *Let Z_1, \dots, Z_k be i.i.d., with $Z_i \in [a, b], i = 1, \dots, k, a < b$ and $\hat{\mu}$ be the empirical estimate $\hat{\mu} = \frac{\sum_{i=1}^k Z_i}{k}$ of $\mathbb{E}[Z] = \mathbb{E}[Z_i]$. Then:*

$$\mathbb{P}[\hat{\mu} - \mathbb{E}[Z] \geq \epsilon] \leq \exp\left(-\frac{2k\epsilon^2}{(b-a)^2}\right) \quad (13)$$

and

$$\mathbb{P}[\hat{\mu} - \mathbb{E}[Z] \leq -\epsilon] \leq \exp\left(-\frac{2k\epsilon^2}{(b-a)^2}\right) \quad (14)$$

hold for all $\epsilon \geq 0$.

In our case we have $k = m$, $Z_i = \mathbb{P}[\hat{Y} = Y_i | \mathcal{C}_\alpha(X_i), Y_i] \in (0, 1)$. Hence, for the empirical estimate $\hat{\mu} = \hat{\mu}_\alpha$ of $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha]$ and its error $\epsilon = \epsilon_{\alpha, \delta}$:

$$\mathbb{P}[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \geq \epsilon_{\alpha, \delta}] \leq \exp\left(-\frac{2m\epsilon_{\alpha, \delta}^2}{(1-0)^2}\right) \quad (15)$$

and

$$\mathbb{P}[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \leq -\epsilon_{\alpha, \delta}] \leq \exp\left(-\frac{2m\epsilon_{\alpha, \delta}^2}{(1-0)^2}\right) \quad (16)$$

hold. Further, if we set

$$\delta = \exp(-2m\epsilon_{\alpha, \delta}^2), \quad (17)$$

then

$$1 - \mathbb{P} \left[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \leq \epsilon_{\alpha, \delta} \right] \leq \delta \Rightarrow \mathbb{P} \left[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \leq \epsilon_{\alpha, \delta} \right] \geq 1 - \delta \quad (18)$$

and

$$1 - \mathbb{P} \left[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \geq -\epsilon_{\alpha, \delta} \right] \leq \delta \Rightarrow \mathbb{P} \left[\hat{\mu}_\alpha - \mathbb{P}[\hat{Y} = Y; \mathcal{C}_\alpha] \geq -\epsilon_{\alpha, \delta} \right] \geq 1 - \delta \quad (19)$$

hold for any $\epsilon_{\alpha, \delta} \geq 0$. As follows, based on Eq. 17:

$$\delta = \exp(-2m\epsilon_{\alpha, \delta}^2) \Rightarrow \log \frac{1}{\delta} = 2m\epsilon_{\alpha, \delta}^2 \Rightarrow \epsilon_{\alpha, \delta}^2 = \frac{\log \frac{1}{\delta}}{2m} \Rightarrow \epsilon_{\alpha, \delta} = \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

C Proof of Theorem 2

Let $\mathcal{C}_{\alpha_1}, \mathcal{C}_{\alpha_2}$ be two set-valued predictors constructed with standard conformal prediction with coverage $1 - \alpha_1$ and $1 - \alpha_2$ respectively. Note that they depend on the calibration data \mathcal{D}_{cal} and thus are random variables. Based on Theorem 1, it holds that:

$$1 - \alpha_1 \leq \mathbb{P}[Y \in \mathcal{C}_{\alpha_1}(X)] \leq 1 - \alpha_1 + \frac{1}{m+1} \quad (20)$$

and

$$1 - \alpha_2 \leq \mathbb{P}[Y \in \mathcal{C}_{\alpha_2}(X)] \leq 1 - \alpha_2 + \frac{1}{m+1}, \quad (21)$$

where the randomness is in the calibration data \mathcal{D}_{cal} and the new sample they help predict.

Since $\alpha_1 < \alpha_2$, we have that $\hat{q}_{\alpha_2} = \frac{\lceil (m+1)(1-\alpha_2) \rceil}{m} < \frac{\lceil (m+1)(1-\alpha_1) \rceil}{m} = \hat{q}_{\alpha_1}$. We know that $\mathcal{C}_{\alpha_2}(X) = \{y : s(X, y) < \hat{q}_{\alpha_2}\}$ and $\mathcal{C}_{\alpha_1}(X) = \{y : s(X, y) < \hat{q}_{\alpha_1}\}$. So for any x , it holds that $\mathcal{C}_{\alpha_2}(x) \subseteq \mathcal{C}_{\alpha_1}(x)$. By the definition of $\mathcal{C}_{\alpha_1, \alpha_2}(X)$ in Eq. 10, we have that:

$$\mathcal{C}_{\alpha_1, \alpha_2}(X) = \{y : \hat{q}_{\alpha_2} < s(X, y) \leq \hat{q}_{\alpha_1}\} = \mathcal{C}_{\alpha_1}(X) \setminus \mathcal{C}_{\alpha_2}(X),$$

and therefore

$$\mathbb{P}[Y \in \mathcal{C}_{\alpha_1, \alpha_2}(X)] = \mathbb{P}[Y \in \mathcal{C}_{\alpha_1}(X)] - \mathbb{P}[Y \in \mathcal{C}_{\alpha_2}(X)]. \quad (22)$$

Combining the upper-bound in Eq. 20 and the lower-bound in Eq. 21, we have¹⁰

$$\mathbb{P}[Y \in \mathcal{C}_{\alpha_1}(X)] - \mathbb{P}[Y \in \mathcal{C}_{\alpha_2}(X)] \leq \alpha_2 - \alpha_1 + \frac{1}{m+1}. \quad (23)$$

Similarly, combining the lower-bound in Eq. 20 and the upper-bound in Eq. 21, we have

$$\alpha_2 - \alpha_1 - \frac{1}{m+1} \leq \mathbb{P}[Y \in \mathcal{C}_{\alpha_1}(X)] - \mathbb{P}[Y \in \mathcal{C}_{\alpha_2}(X)]. \quad (24)$$

From Eqs. 22, 23 and 24, we have that:

$$\alpha_2 - \alpha_1 - \frac{1}{m+1} \leq \mathbb{P}[Y \in \mathcal{C}_{\alpha_1, \alpha_2}(X)] \leq \alpha_2 - \alpha_1 + \frac{1}{m+1}.$$

¹⁰We use the equality in Eq.22 instead of applying a union bound to derive the results in Eq. 23 and Eq. 24 when combining the bounds.

D Implementation Details

To implement our algorithms and run all the experiments on synthetic and real data, we used PyTorch 1.8.1, NumPy 1.20.1 and Scikit-learn 1.0.2 on Python 3.7.3. For reproducibility, we use a fixed random seed in all random procedures. Moreover, we set $\delta = 0.1$ everywhere.

Synthetic prediction tasks. We create $4 \times 3 = 12$ different prediction tasks, where we vary the number of labels $n \in \{10, 50, 100\}$ and the level of difficulty of the task. More specifically, for each value of n , we create four different tasks of increasing difficulty where the success probability of the logistic regression classifier is $\mathbb{P}[Y' = Y] = 0.9, 0.7, 0.5$ and 0.3 , respectively.

To create each task, we use the function `make_classification` of the Scikit-learn library. This function allows the creation of data for synthetic prediction tasks with very particular user-defined characteristics, through the generation of clusters of normally distributed points on the vertices of a multidimensional hypercube. The number of the dimensions of the hypercube indicates the number of informative features of each sample, which in our case we set at 15 for all prediction tasks. Linear combinations of points, *i.e.*, the informative features, are used to create redundant features, the number of which we set at 5. The difficulty of the prediction task is controlled through the size of the hypercube, with a multiplicative factor, namely `clas_sep`, which we tuned accordingly for each value n so that the success probability of the logistic regression classifier above spans a wide range of values across tasks. All the selected values of this parameter can be found in the configuration file `config.py` in the code. Finally, we set the proportion of the samples assigned to each label, *i.e.*, the function parameter `weights`, using a Dirichlet distribution of order n with parameters $\alpha_1 = \dots = \alpha_n = 1$.

E Additional Experiments on Synthetic Data

To complement the results in Table 1 in the main paper, we experiment with additional prediction tasks with different number of labels n and amount of calibration and estimation data m . For each value of n and m , we estimate the average success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}_1, \hat{\alpha}_2}]$ during test for four different experts using our system, each with a different success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, on four prediction tasks where the classifier achieves a different success probability $\mathbb{P}[Y' = Y]$. Figure 4 summarizes the results.

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.56	0.72	0.85	0.94
0.5	0.67	0.80	0.89	0.96
0.7	0.79	0.87	0.93	0.98
0.9	0.92	0.95	0.97	0.99

(a) $n = 50, m = 1,200$

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.63	0.76	0.87	0.95
0.5	0.73	0.83	0.91	0.96
0.7	0.82	0.90	0.95	0.98
0.9	0.93	0.95	0.98	0.99

(b) $n = 100, m = 1,200$

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.40	0.57	0.74	0.91
0.5	0.55	0.67	0.80	0.93
0.7	0.72	0.79	0.88	0.96
0.9	0.91	0.92	0.94	0.98

(c) $n = 10, m = 400$

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.56	0.71	0.84	0.94
0.5	0.67	0.80	0.89	0.96
0.7	0.80	0.88	0.94	0.98
0.9	0.91	0.94	0.97	0.99

(d) $n = 50, m = 400$

$\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$	$\mathbb{P}[Y' = Y]$			
	0.3	0.5	0.7	0.9
0.3	0.63	0.76	0.88	0.95
0.5	0.73	0.84	0.92	0.96
0.7	0.83	0.89	0.95	0.98
0.9	0.92	0.96	0.98	0.99

(e) $n = 100, m = 400$

Figure 4: Average success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{C}_{\hat{\alpha}_1, \hat{\alpha}_2}]$ during test for four different experts using our system, each with a different success probability $\mathbb{P}[\hat{Y} = Y; \mathcal{Y}]$, on four prediction tasks where the classifier achieves a different success probability $\mathbb{P}[Y' = Y]$. Each table corresponds to a different number of label values n and calibration and estimation set size m . For readability, each cell shows only the average since the standard errors are all below 10^{-2} .