

PAC CONFIDENCE SETS FOR DEEP NEURAL NETWORKS VIA CALIBRATED PREDICTION

Sangdon Park
University of Pennsylvania
sangdonp@cis.upenn.edu

Osbert Bastani
University of Pennsylvania
obastani@seas.upenn.edu

Nikolai Matni
University of Pennsylvania
nmatni@seas.upenn.edu

Insup Lee
University of Pennsylvania
lee@cis.upenn.edu

ABSTRACT

We propose an algorithm combining calibrated prediction and generalization bounds from learning theory to construct confidence sets for deep neural networks with PAC guarantees—i.e., the confidence set for a given input contains the true label with high probability. We demonstrate how our approach can be used to construct PAC confidence sets on ResNet for ImageNet, a visual object tracking model, and a dynamics model for the half-cheetah reinforcement learning problem.¹

1 INTRODUCTION

A key challenge facing deep neural networks is that they do not produce reliable confidence estimates, which are important for applications such as safe reinforcement learning (Berkenkamp et al., 2017), guided exploration (Malik et al., 2019), and active learning (Gal et al., 2017).

We consider the setting where the test data follows the same distribution as the training data (i.e., we do *not* consider adversarial examples designed to fool the network (Szegedy et al., 2014)); even in this setting, confidence estimates produced by deep neural networks are notoriously unreliable (Guo et al., 2017). One intuition for this shortcoming is that unlike traditional supervised learning algorithms, deep learning models typically overfit the training data (Zhang et al., 2017). As a consequence, the confidence estimates of deep neural networks are flawed even for test data from the training distribution since, by construction, they overestimate the likelihood of the training data.

A promising approach to addressing this challenge is *temperature scaling* (Platt, 1999). This approach takes as input a trained neural network $f_{\hat{\phi}}(y | x)$ —i.e., whose parameters $\hat{\phi}$ have already been fit to a training dataset Z_{train} —which produces unreliable probabilities $f_{\hat{\phi}}(y | x)$. Then, this approach rescales these confidence estimates based on a validation dataset to improve their “calibration”. More precisely, this approach fits confidence estimates of the form

$$f_{\hat{\phi},\tau}(y | x) \propto \exp(\tau \log f_{\hat{\phi}}(y | x)),$$

where $\tau \in \mathbb{R}_{>0}$ is a *temperature scaling* parameter that is fit based on the validation dataset. The goal is to choose τ to minimize *calibration error*, which roughly speaking measures the degree to which the reported error rate differs from the actual error rate.

The key insight is that in the temperature scaling approach, only a single parameter τ is fit to the validation data—thus, unlike fitting the original neural network, the temperature scaling algorithm comes with generalization guarantees based on traditional statistical learning theory.

Despite the improved generalization guarantees, these confidence estimates still do not come with theoretical guarantees. We are interested in producing *confidence sets* that satisfy statistical guarantees while being as small as possible. Given a test input $x \in \mathcal{X}$, a confidence set $C_T(x) \subseteq \mathcal{Y}$

¹Our code is available at <https://github.com/sangdon/PAC-confidence-set>.

	$ C(x) = 1$	$5 \leq C(x) \leq 10$	$50 \leq C(x) \leq 100$	$ C(x) \geq 200$
airship				
zebra				

Table 1: ImageNet images with varying ResNet confidence set sizes. The confidence set sizes are on the top. The true label is on the left-hand side. Incorrectly labeled images are boxed in red.

(parameterized by $\epsilon \in (0, 1]$) should contain the true label for at least a $1 - \epsilon$ fraction of cases:

$$P_{(x,y) \sim D} [y \in C_T(x)] \geq 1 - \epsilon.$$

Since we are fitting a parameter θ based on Z_{val} , we additionally incur a probability of failure due to the randomness Z_{val} . In other words, given $\epsilon \in (0, 1]$, we aim to obtain probably approximately correct (PAC) confidence sets $C_T(x) \subseteq Y$ satisfying the guarantee

$$P_{Z_{\text{val}} \sim D^n} P_{(x,y) \sim D} (y \in C_T(x)) \geq 1 - \epsilon.$$

Indeed, techniques from statistical learning theory (Vapnik, 1999) can be used to do so (Vovk, 2013).

There are a number of reasons why confidence sets can be useful. First, they can be used to inform safety critical decision making. For example, consider a doctor who uses prediction tools to help perform diagnosis. Having a confidence set would both help the doctor estimate the confidence of the prediction (i.e., smaller confidence sets imply higher confidence), but also give a sense of the set of possible diagnoses. Second, having a confidence set can be useful for reasoning about safety since they contain the true outcome with high probability. For instance, robots may use a confidence set over predicted trajectories to determine whether it is safe to act with high probability. As a concrete example, consider a self-driving car that uses a deep neural network to predict the path that a pedestrian might take. We require that the self-driving car avoid the pedestrian with high probability, which it can do by avoiding all possible paths in the predicted confidence set.

Contributions. We propose an algorithm combining calibrated prediction and statistical learning theory to construct PAC confidence sets for deep neural networks (Section 3). We propose instantiations of this framework in the settings of classification, regression, and learning models for reinforcement learning (Section 3.6). Finally, we evaluate our approach on three benchmarks: ResNet (He et al., 2016) for ImageNet (Russakovsky et al., 2015), a model (Held et al., 2016) learned for a visual object tracking benchmark (Wu et al., 2013), and a probabilistic dynamics model (Chua et al., 2018) learned for the half-cheetah environment (Brockman et al., 2016) (Section 4). Examples of ImageNet images with different sized ResNet confidence sets are shown in Table 1. As can be seen, our confidence sets become larger and the images become more challenging to classify. In addition, we show predicted confidence sets for ResNet in Table 2, as well as predicted confidence sets for the visual object tracking model in Table 3.

Related work. There has been work on constructing confidence sets with theoretical guarantees. Oftentimes, these guarantees are asymptotic rather than finite sample (Steinberger & Leeb, 2016; 2018). Alternatively, there has been work focused on predicting confidence sets with a given expected size (Denis & Hebiri, 2017).

More relatedly, there has been recent work on obtaining PAC guarantees. For example, there has been some work specific prediction tasks such as binary classification (Lei, 2014; Wang & Qiao,




$1 \leq C(x) < 5$	$5 \leq C(x) < 10$	$10 \leq C(x) < 20$
 <p>king penguin</p>	 <p>shopping basket</p>	 <p>chambered nautilus</p>
	<p>Chihuahua toy terrier Italian greyhound Boston bull miniature pinscher</p> <p>English springer Welsh springer spaniel collie; boxer Saint Bernard Leonberg</p> <p>face powder hamper lotion; packet shopping basket</p>	<p>banded gecko common iguana American chameleon whiptail; agama frilled lizard; alligator lizard green lizard African chameleon Komodo dragon</p> <p>altar; analog clock bell cote castle church cinema dome monastery palace vault wall clock</p> <p>barber chair hand blower medicine chest paper towel plunger shower curtain soap dispenser toilet seat tub; washbasin washer toilet tissue</p>

Table 2: Confidence sets of ImageNet images with varying ResNet confidence set sizes. The predicted confidence set is shown to the right of the corresponding input image. The true label is shown in red, and the predicted label is shown with a hat. See Table 5 in Appendix D for more examples.

Table 3: Visualization of confidence sets for the tracking dataset (Wu et al., 2013), including the ground truth bounding box (white), the bounding box predicted by the original neural network (Held et al., 2016) (red), and the bounding box produced using our confidence set predictor (green). We have overapproximated the predicted ellipsoid confidence set with a box. Our bounding box contains the ground truth bounding box with high probability. See Table 9 in Appendix D for more examples.

2018). There has also been work in the setting of regression (Lei et al., 2018; Barber et al., 2019). However, in this case, the confidence sets are fixed in size—i.e., they do not depend on the input x (Barber et al., 2019). Furthermore, they make stability assumptions about the learning algorithm (though they achieved improved rates by doing so) (Lei et al., 2018; Barber et al., 2019).

The most closely related work is conformal prediction (Papadopoulos, 2008; Vovk, 2013). Like our approach, this line of work provides a way to construct confidence sets from a given confidence predictor, and provided PAC guarantees for the validity of these confidence sets. Indeed, with some work, our generalization bound Theorem 1 can be shown to be equivalent to Theorem 1 in Vovk (2013). In contrast to their approach, we proposed to use calibrated prediction to construct confidence predictors that can suitably be used with deep neural networks. Furthermore, our approach

makes explicit the connections to temperature scaling and as well as to generalization bounds from statistical learning theory (Vapnik, 1999). In addition, unlike our paper, they do not explicitly provide an efficient algorithm for constructing confidence sets. Finally, we also propose an extension to the case of learning models for reinforcement learning.

Finally, we build on a long line of work on calibrated prediction which aims to construct “calibrated” probabilities (Murphy, 1972; DeGroot & Fienberg, 1983; Platt, 1999; Zadrozny & Elkan, 2001; 2002; Naeini et al., 2015; Kuleshov & Liang, 2015). Roughly speaking, probabilities are calibrated if events happen at rates equal to the predicted probabilities. This work has recently been applied to obtaining confidence estimates for deep neural networks (Guo et al., 2017; Kuleshov et al., 2018; Pearce et al., 2018), including for learned models for reinforcement learning (Malik et al., 2019). However, these approaches do not come with PAC guarantees.

2 PAC CONFIDENCE SETS

Our goal is to estimate confidence sets that are as small as possible, while simultaneously ensuring that they are probably approximately correct (PAC) (Valiant, 1984). Essentially, a confidence set is correct if it contains the true label. More precisely, let \mathcal{X} be the inputs and \mathcal{Y} be the labels, and let D be a distribution over $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. A confidence set predictor is a function $C : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ such that $C(x) \subseteq \mathcal{Y}$ is a set of labels; we denote the set of all confidence set predictors by \mathcal{C} . For a given example $(x; y) \in D$, we say C is correct if $y \in C(x)$. Then, the error of C is

$$L(C) = \mathbb{P}_{(x,y) \sim D} [y \notin C(x)] \quad (1)$$

Finally, consider an algorithm A that takes as input a validation set $Z_{\text{val}} \subseteq \mathcal{Z}$ consisting of n i.i.d. samples $(x; y) \in D$, and outputs a confidence set predictor \hat{C} . Given $\epsilon \in (0, 1]$ and $\delta \in (0, 1]$, we say that A is probably approximately correct (PAC) if

$$\mathbb{P}_{Z_{\text{val}} \sim D^n} [L(\hat{C}) > \epsilon] < \delta \quad \text{where } \hat{C} = A(Z_{\text{val}}) \quad (2)$$

Our goal is to design an algorithm A that satisfies (2) while constructing confidence sets $C(x)$ that are as “small in size” as possible on average. The size of $C(x)$ depends on the domain. For classification, we consider confidence sets that are arbitrary subsets of labels $\mathcal{Y} = \{1; \dots; Y\}$, and we measure the size by $|C(x)| \in \mathbb{N}$ —i.e., the number of labels in $C(x)$. For regression, we consider confidence sets that are intervals $C(x) = [a; b] \subseteq \mathcal{Y} = \mathbb{R}$, and we measure size by $b - a$ —i.e., the length of the predicted interval. Note that there is an intrinsic tradeoff between satisfying (2) and average size of $C(x)$ —larger confidence sets are more likely to satisfy (2).

3 PAC ALGORITHM FOR CONFIDENCE SET CONSTRUCTION

Our algorithm is formulated in the empirical risk framework. Typically, this framework refers to empirical risk minimization. In our setting, such an algorithm would take as input (i) a parametric family of confidence set predictors $\mathcal{C} = \{C_\theta \mid \theta \in \Theta\}$, where Θ is the parameter space, and (ii) a training set $Z_{\text{val}} \subseteq \mathcal{Z}$ of n i.i.d. samples $(x; y) \in D$, and output the confidence set predictor \hat{C} , where \hat{C} minimizes the empirical risk:

$$\hat{C} = \arg \min_{\mathcal{C}} \hat{L}(C; Z_{\text{val}}) \quad \text{where} \quad \hat{L}(C; Z_{\text{val}}) = \frac{1}{n} \sum_{(x,y) \in Z_{\text{val}}} \mathbb{1}[y \notin C(x)]$$

Here, $\mathbb{1}[\cdot] \in \{0, 1\}$ is the indicator function, and the empirical risk \hat{L} is an estimate of the confidence set error (1) based on the validation set Z_{val} .

However, our algorithm does not minimize the empirical risk. Rather, recall that our goal is to minimize the size of the predicted confidence sets given a PAC constraint on the true risk based on the given PAC parameters $\epsilon \in (0, 1]$ and $\delta \in (0, 1]$ and the number of available validation samples $n = |Z_{\text{val}}|$. Thus, the risk shows up as a constraint in the optimization problem, and the objective is instead to minimize the size of the predicted confidence sets:

$$\hat{C} = \arg \min_{\mathcal{C}} S(\hat{C}) \quad \text{subj. to} \quad \hat{L}(C; Z_{\text{val}}) \leq \epsilon \quad (3)$$

At a high level, the value $\beta = (n; \epsilon) \in \mathbb{R}_{>0}$ is chosen to enforce the PAC constraint, and is based on generalization bounds from statistical learning theory (Valiant, 1984). Furthermore, following the temperature scaling approach (Platt, 1999), the parameter space is chosen to be as small as possible (in particular, one dimensional) to enable good generalization. Finally, our choice of size metric S follows straightforwardly based on our choice of parameter space. In the remainder of this section, we describe the choices of (i) parameter space Θ , and (iii) confidence level $(n; \epsilon)$ in more detail, as well as how to solve (3) given these choices.

3.1 CHOICE OF PARAMETER SPACE

Probability forecasters. Our construction of the parametric family of confidence set predictors assumes given a probability forecaster $f : X \rightarrow \mathcal{P}_Y$, where \mathcal{P}_Y is a space of probability distributions over Y . Given such an f , we use $f(y|x)$ to denote the probability of label y under distribution $f(x)$. Intuitively, $f(y|x)$ should be the probability (or probability density) that the true label for a given input—i.e., $f(y|x) = \mathbb{P}_{(x;Y) \sim D}[Y = y | X = x]$. For example, in classification, we can choose \mathcal{P}_Y to be the space of categorical distributions over \mathcal{Y} and f may be a neural network whose last layer is a softmax layer with $|\mathcal{Y}|$ outputs. Then $f(y|x) = f(x)_y$. Alternatively, in regression, we can choose \mathcal{P}_Y to be the space of Gaussian distributions, and f may be a neural network whose last layer outputs the values $(\mu; \sigma^2) \in \mathbb{R}_{>0}$ of a Gaussian distribution. Then, $f(y|x) = N(x; (\mu; \sigma^2))$, where $(\mu; \sigma^2) = f(x)$, and $N(\cdot; (\mu; \sigma^2))$ is the Gaussian density function with mean μ and variance σ^2 .

Training a probability forecaster. To train a probability forecaster, we use a standard approach to calibrated prediction that combines maximum likelihood estimation with temperature scaling². First, we consider a parametric model family $\mathcal{F} = \{f_\theta | \theta \in \Theta\}$, where Θ is the parameter space. Note that Θ can be high-dimensional—e.g., the weights of a neural network model. Given a training set $Z_{\text{train}} = Z$ of m i.i.d. samples $(x; y) \in D$, the maximum likelihood estimate (MLE) of θ is

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \ell(\theta; Z_{\text{train}}) \quad \text{where} \quad \ell(\theta; Z_{\text{train}}) = \sum_{(x; y) \in Z_{\text{train}}} \log f_\theta(y|x) \quad (4)$$

We could now use $\hat{\theta}$ as the probability forecaster. However, the problem with directly using $\hat{\theta}$ is that because $\hat{\theta}$ may be high-dimensional, it often overfits the training data Z_{train} . Thus, the probabilities are typically overconfident compared to what they should be.

To reduce their confidence, we use the temperature scaling approach to calibrate the predicted probabilities (Platt, 1999; Guo et al., 2017). Intuitively, this approach is to train an MLE estimate using exactly the same approach used to train $\hat{\theta}$ but using a single new parameter $\beta \in \mathbb{R}_{>0}$. The key idea is that this time, the model family is based on the parameter $\hat{\theta}$ from (4). In other words, the “shape” of the probabilities forecast by $\hat{\theta}$ are preserved, but their exact values are shifted.

More precisely, consider the model family $\mathcal{F}^\beta = \{f_{\hat{\theta}; \beta} | \beta \in \mathbb{R}_{>0}\}$, where

$$f_{\hat{\theta}; \beta}(y|x) = \exp(\beta \log f_{\hat{\theta}}(y|x)) / \sum_{y' \in \mathcal{Y}} \exp(\beta \log f_{\hat{\theta}}(y'|x))$$

Then, we have the following MLE for:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}_{>0}} \ell^0(\beta; Z_{\text{train}}^0) \quad \text{where} \quad \ell^0(\beta; Z_{\text{train}}^0) = \sum_{(x; y) \in Z_{\text{train}}^0} \log f_{\hat{\theta}; \beta}(y|x) \quad (5)$$

Note that $\hat{\beta}$ is estimated based on a second training set Z_{train}^0 . Because we are only fitting a single parameter, this training set can be much smaller than the training set used to train $\hat{\theta}$.

Parametric family of confidence set predictors. Finally, given a probability forecaster f , we consider one dimensional parameter space \mathbb{R} ; in an analogy to the temperature scaling technique for calibrated prediction, we denote this parameter by β . In particular, we assume a confidence probability predictor f is given, and consider

$$C_T(x) = \{y \in \mathcal{Y} | f(y|x) \geq e^{-T} g\}$$

²A priori, it is not obvious that using temperature scaling can improve our confidence set predictor; we give a detailed discussion in Appendix A.1.

In other words, $\mathcal{C}_T(x)$ is the set of y with high probability given x according to \mathbb{P} . Considering this scalar parameter space, we denote the minimum of (3) by

3.2 CHOICE OF SIZE METRIC $S(T)$

To choose the size metric $S(T)$, we note that for our chosen parametric family of confidence set predictors, smaller values correspond to uniformly smaller confidence sets—i.e.,

$$T \leq T' \implies \mathcal{C}_T(x) \subseteq \mathcal{C}_{T'}(x).$$

Thus, we can simply choose the size metric to be

$$S(T) = T. \quad (6)$$

This choice minimizes the size of the confidence sets produced by our algorithm.

3.3 CHOICE OF CONFIDENCE LEVEL $(\epsilon; \delta)$

Naive approach based on VC generalization bound A naive approach to choosing $(\epsilon; \delta)$ is to do so based on the VC dimension generalization bound (Vapnik, 1999). It is not hard to show that the problem of estimating \mathbb{P} is equivalent to a binary classification problem, and that the VC dimension of \mathcal{C} for this problem is d . Thus, the VC dimension bound implies that for $\forall \epsilon \in [0, 1]$,

$$\mathbb{P}_{Z_{\text{val}} \sim D^n} L(\hat{C}_T) \leq L(C_T; Z_{\text{val}}) + \frac{\log(2n) + 1}{n} \log\left(\frac{1}{\epsilon}\right). \quad (7)$$

The details of this equivalence are given in Appendix B.2. Then, suppose we choose

$$(\epsilon; \delta) = \left(\frac{\log(2n) + 1}{n} \log\left(\frac{1}{\delta}\right); \delta\right).$$

With this choice, for the solution \hat{T} of (3) with $(\epsilon; \delta) = \left(\frac{\log(2n) + 1}{n} \log\left(\frac{1}{\delta}\right); \delta\right)$, the constraint in (3) ensures that $L(\hat{C}_{\hat{T}}; Z_{\text{val}}) \leq \delta$. Together with the VC generalization bound (7), we have

$$\mathbb{P}_{Z_{\text{val}} \sim D^n} L(\hat{C}_{\hat{T}}) < \epsilon;$$

which is exactly desired the PAC constraint on our predicted confidence sets.

Direct generalization bound. In fact, we can get better choices of $(\epsilon; \delta)$ by directly bounding generalization error. For instance, in the realizable setting (i.e., we always have $L(C_{\mathbb{P}}; Z_{\text{val}}) = 0$), we can get rates of $\epsilon = \mathcal{O}(1/n)$ instead of $\epsilon = \mathcal{O}(1/n^2)$ (Kearns & Vazirani, 1994); see Appendix A.2 for details. We can achieve these rates by choosing $\delta = 0$, but then, the PAC guarantees we obtain may actually be stronger than desired (i.e., $\epsilon < \epsilon_{\text{target}}$). Intuitively, we can directly prove a bound that interpolates between the realizable setting and the VC generalization bound—in particular:

Theorem 1. For any $\epsilon \in [0, 1]$, $n \geq N_{\epsilon}$, and $k \in [0, 1]$, we have

$$\mathbb{P}_{Z_{\text{val}} \sim D^n} L(\hat{C}_{\hat{T}}) < \epsilon + \sum_{i=0}^{k-1} \frac{n^i}{i!} (1 - \epsilon)^{n-i};$$

where \hat{T} is the solution to (3) with $(\epsilon; \delta) = (k\epsilon; 0)$.³

We give a proof in Appendix B.2. Based on Theorem 1, we can choose

$$(\epsilon; \delta) = \max_{k \in [0, 1]} \{k\epsilon \mid \sum_{i=0}^{k-1} \frac{n^i}{i!} (1 - \epsilon)^{n-i} < \epsilon\}. \quad (8)$$

3.4 THEORETICAL GUARANTEES

We have the following guarantee, which follows straightforwardly from Theorem 1:

Corollary 1. Let \hat{T} be the solution to (3) for $(\epsilon; \delta) = (\epsilon; 0)$ chosen according to (8). Then, we have

$$\mathbb{P}_{Z_{\text{val}} \sim D^n} L(\hat{C}_{\hat{T}}) < \epsilon.$$

In other words, our algorithm is probably approximately correct.

³The theorem statement relies on additional standard technical conditions; see Appendix B.1.

Algorithm 1 Algorithm for solving (3).

```

procedure ESTIMATECONFIDENCESETPREDICTOR( $Z_{\text{train}}; Z_{\text{train}}^0; Z_{\text{val}}$ )
  Estimate  $\hat{f}, \hat{\lambda}$  using (4) and (5), respectively
  Compute  $(n; \tau)$  according to (8) by enumerating  $\tau \in \{0, 1, \dots, n\}$ 
  Let  $k = n - (n; \tau)$  (note that  $\tau \in \{0, 1, \dots, n\}$ )
  Sort  $(x; y) \in Z_{\text{val}}$  in ascending order of  $\hat{f}_{\hat{\lambda}}(y | x)$ 
  Let  $(x_{k+1}; y_{k+1})$  be the  $(k+1)$ st element in the sorted  $Z_{\text{val}}$ 
  Solve (3) by choosing  $\hat{\tau} = \log \hat{f}_{\hat{\lambda}}(y_{k+1} | x_{k+1})$ 
  Return  $C_{\hat{\tau}} = \{x \in Z_{\text{val}} \mid \hat{f}_{\hat{\lambda}}(y | x) < e^{\hat{\tau}}\}$ 
end procedure

```

3.5 PRACTICAL IMPLEMENTATION

Our algorithm for estimating a confidence set predictor is summarized in Algorithm 1. The algorithm solves the optimization problem (3) using the choices $\hat{f}(T)$, and $(n; \tau)$ described in the preceding sections. There are two key implementation details that we describe here.

Computing $(n; \tau)$. To compute $(n; \tau)$, we need to solve (8). A straightforward approach is to enumerate all possible choices $\tau \in \{0, 1, \dots, n\}$. There are two optimizations. First, the objective is monotone increasing in τ , so we can enumerate in ascending order until the constraint no longer holds. Second, rather than re-compute the left-hand side of the constraint $\sum_{(x,y) \in Z_{\text{val}}} \mathbb{1}_{\hat{f}_{\hat{\lambda}}(y|x) < e^{\tau}}$, we can accumulate the sum as we iterate over τ . We can also incrementally compute $\hat{f}, \hat{\lambda}$, and $(1 - \epsilon)^{n-i}$. For numerical stability, we perform these computations in log space.

Solving (3). To solve (3), note that the constraint in (3) is equivalent to

$$\sum_{(x,y) \in Z_{\text{val}}} \mathbb{1}_{\hat{f}_{\hat{\lambda}}(y|x) < e^{\tau}} = n - (n; \tau) \quad \text{where } E(x; y; T) = \mathbb{1}_{\hat{f}_{\hat{\lambda}}(y|x) < e^{\tau}} : \quad (9)$$

Also, note that $k = n - (n; \tau)$ is an integer due to the definition of $(n; \tau)$ in (8). Thus, we can interpret (9) as saying that $E(x; y; T) = 1$ for at most k of the points $(x; y) \in Z_{\text{val}}$.

In addition, note that $E(x; y; T)$ decreases monotonically as $\hat{f}_{\hat{\lambda}}(y | x)$ becomes larger. Thus, we can sort the points $(x; y) \in Z_{\text{val}}$ in ascending order of $\hat{f}_{\hat{\lambda}}(y | x)$, and require that only the k smallest points $(x; y)$ in this list satisfy $E(x; y; T) = 1$. In particular, letting $(x_{k+1}; y_{k+1})$ be the $(k+1)$ st point, (9) is equivalent to

$$\hat{f}_{\hat{\lambda}}(y_{k+1} | x_{k+1}) < e^{\tau} : \quad (10)$$

In other words, this constraint says that $(x_{k+1}; y_{k+1}) \in C_T(x_{k+1})$. Finally, the solution $\hat{\tau}$ to (3) is the smallest τ that satisfies (10), which is the τ that makes (10) hold with equality—i.e.,

$$\hat{\tau} = \log \hat{f}_{\hat{\lambda}}(y_{k+1} | x_{k+1}) : \quad (11)$$

We have assumed $\hat{f}_{\hat{\lambda}}(y_{k+1} | x_{k+1}) > \hat{f}_{\hat{\lambda}}(y_k | x_k)$; if not, we decrement k until this holds.

3.6 PROBABILITY FORECASTERS FOR SPECIFIC TASKS

We briefly discuss the architectures we use for probability forecasters for various tasks. We give details, including how we measure the sizes of predicted confidence sets, in Appendix C. We consider three tasks: classification, regression, and model-based reinforcement learning. For classification, we use the standard approach of using a soft-max layer to predict label probabilities $f(y | x)$. For regression, we also use a standard approach where the neural network predicts both the mean $\mu(x)$ and covariance $\Sigma(x)$ of a Gaussian distribution $\mathcal{N}(x; (\mu(x); \Sigma(x)))$; then, $f(y | x) = \mathcal{N}(y; (\mu(x); \Sigma(x)))$ is the probability density of according to this Gaussian distribution.

Finally, for model-based reinforcement learning, our goal is to construct confidence sets over trajectories predicted using a learned model of the dynamics. We consider unknown dynamics



Figure 1: Results on ResNet for ImageNet with $n = 20000$. Default parameters are $\epsilon = 0.01$ and $\delta = 10^{-5}$. We plot the median and min/max confidence set sizes. (a) Ablation study; “calibrated predictor” (i.e., use $\hat{\lambda}_\lambda$ instead of λ), and D is “direct bound” (i.e., use Theorem 1 instead of the VC generalization bound). (b) Restricted to correctly vs. incorrectly labeled images. (c) Varying ϵ . (d) Varying δ .

$g(x^0_j x; u)$ mapping a state-action pair $(x; u)$ to a distribution over states \mathcal{X} , and consider a known (and fixed) policy $(u_j x)$ mapping a given state to a distribution over actions $\mathcal{U} \subseteq \mathbb{R}^{d_u}$. Then, we let $f(x^0_j x) = E_{(u_j x)}[g(x^0_j x; u)]$ denote the (unknown) closed-loop dynamics.

Next, we consider a forecast $f(x^0_j x)$ of the form $f(x^0_j x) = N(x^0; \mu(x); \Sigma(x))$, and our goal is to construct confidence sets for the predictions of f . However, we want to do so for not just for one-step predictions, but for predictions over a time horizon N . In particular, given initial state $x_0 \in \mathcal{X}$, we can sample $x_{1:H} = (x_1; \dots; x_H)$ by letting $x_0 = x_0$ and sequentially sampling $x_{t+1} \sim f(\cdot_j x_t)$ for each $t \in \{0; 1; \dots; H-1\}$. Then, our goal is to construct a confidence set that contains $x_{1:H} \in \mathcal{X}^H$ with high probability (over both the randomness in an initial state distribution $x_0 \sim d_0$ and the randomness in f).

To do so, we construct and use a forecaster $f_{1:H}(x_{1:H} | x_0)$ based on f . In principle, this task is a special case of multivariate regression, where the inputs x (i.e., the initial state x_0) and the outputs are $Y = X^H$ (i.e., a predicted trajectory $x_{1:H}$). However, the variance $\Sigma(x)$ predicted by our probability forecaster is only for a single step, and does not take into account the fact that itself is uncertain. Thus, we use a simple heuristic where we accumulate variances over time. More precisely, we construct (i) the predicted mean $\mu_{1:H} = (\mu_1; \dots; \mu_H)$ by $x_0 = x_0$ and $x_{t+1} = f(x_t)$ for $t \in \{0; 1; \dots; H-1\}$, and (ii) the predicted variances $\Sigma_{1:H} = (\Sigma_1; \dots; \Sigma_H)$ by

$$\tilde{\Sigma}_t = \Sigma(x_0) + \Sigma(x_1) + \dots + \Sigma(x_{t-1});$$

We use a probability forecaster $f_{1:H}(x_{1:H} | x_0) = N(x_{1:H}; \mu_{1:H}; \tilde{\Sigma}_{1:H})$ to construct confidence sets.

4 EXPERIMENTS

We describe our experiments on ImageNet (a classification task), a visual object tracking benchmark (a regression task), and the half-cheetah environment (a model-based reinforcement learning task). We give additional results in Appendix D.

(a) (b) (c)

Figure 2: Confidence set sizes for an object tracking benchmark (Wu et al., 2013); we use $n = 5,000$, $\epsilon = 0.01$, and $\delta = 10^{-5}$. (a) Ablation study similar to Figure 3. In (b) and (c), we show how the confidence set sizes produced using our algorithm vary with respect to ϵ , respectively.

ResNet for ImageNet. We use our algorithm to compute confidence sets for ResNet (He et al., 2016) on ImageNet (Russakovsky et al., 2015), for $\epsilon = 0.01$, $\delta = 10^{-5}$, and $n = 20,000$ validation images. We show the results in Figure 1. In (a), we compare our approach to an ablation. In particular, C refers to performing an initial temperature scaling step to calibrate the neural network predictor (i.e., using \hat{f}_λ instead of $\hat{f}_{\hat{\lambda}}$), and (ii) D refers to using Theorem 1 instead of the VC generalization bound. Thus C + D refers to our approach. As can be seen, using calibrated predictor produces a noticeable reduction in the maximum confidence set size.

We also compared to the ablation—i.e., using the VC generalization bound. However, we were unable to obtain valid confidence sets for our choice of ϵ and δ —i.e., (3) is infeasible. That is, using Theorem 1 outperforms using the VC generalization bound since the VC bound is too loose to satisfy the PAC criterion for our choice of parameters. In addition, in Table 6 in Appendix D, we show results for larger choices of ϵ and δ ; these results show that our approach substantially outperforms the ablation based on the VC bound even when the VC bound produces valid confidence sets.

In (b), we show the confidence set sizes for images correctly vs. incorrectly labeled by ResNet. As expected, the sizes are substantially larger for incorrectly labeled images. Finally, in (c) and (d), we show how the sizes vary with ϵ and δ , respectively. As expected, the dependence is much more pronounced (note that is log-scale).

Visual object tracking. We apply our confidence set prediction algorithm to a 2D visual single-object tracking task, which is a multivariate regression problem. Specifically, the input space consists of the previous image, the previous bounding box \mathcal{R}^2 , and the current image. The output space $\mathcal{Y} = \mathcal{R}^4$ is a current bounding box. We use the regression-based tracker from Held et al. (2016), and retrain the regressor neural network to predict the mean and variance of a Gaussian distribution. More precisely, our object tracking model predicts the mean and variance of each bounding box parameter—i.e., $(x_{\min}, y_{\min}, x_{\max}, y_{\max})$. Given this bounding box forecast \hat{f}_λ , we calibrate and estimate a confidence set predictor as described in Section 3.6.

We use the visual object tracking benchmark from Wu et al. (2013) to train and evaluate our confidence set predictor. This benchmark consists of 99 video sequences labeled with ground truth bounding boxes. We randomly split these sequences to form the training set for calibration, validation set for confidence set estimation, and test set for evaluation. For each sequence, a pair of two adjacent frames constitute a single example. Our training dataset contains 20,882 labeled examples, each consisting of a pair of consecutive images and ground truth bounding boxes. The validation set for confidence set estimation and test set contain 22,761 and 22,761 labeled examples, respectively. Figure 2 shows the sizes of the predicted confidence sets; the sizes are measured as described in Section 3.6 for regression tasks. As for ResNet, we omit results for the VC bound ablation since n is too small to get a bound. The trends are similar to the ones for ResNet.

Half-cheetah. We use our algorithm to compute confidence sets for a probabilistic neural network dynamics model (Chua et al., 2018) for the half-cheetah environment (Brockman et al., 2016), for $\epsilon = 0.01$, $\delta = 10^{-5}$, $H = 20$ time steps, and $n = 5,000$ validation rollouts. When using temperature scaling to calibrate \hat{f}_λ to obtain $\hat{f}_{\hat{\lambda}}$, we calibrate each dimension of time steps independently (i.e., we tune H parameters, where H is time horizon). We show the results in Figure 3.

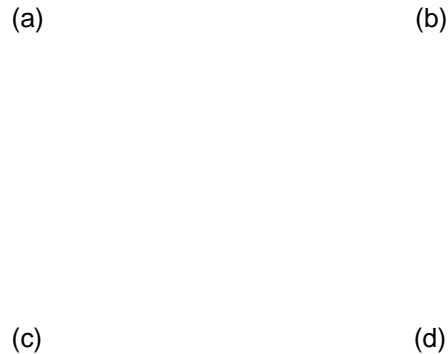


Figure 3: Results on the dynamics model for the half-cheetah with 5000. Default parameters are $\alpha = 0.01$ and $\beta = 10^{-5}$. (a) Ablation study: A is “accumulated variance” (i.e., for each $t \in \{1, \dots, 20\}$, use \tilde{x}_t instead of $x_t = (x_{t-1})$), and C and D are as for ResNet. We plot the median and min/max confidence set sizes (see Section 3.6), averaged across $t \in \{1, \dots, 20\}$. (b) Same ablations, but with per time step size. We plot the average size of the confidence set for the predicted state x_t on step t , as a function of $t \in \{1, \dots, 20\}$. (c) Varying α , and (d) varying β .

In (a), we compare to two ablations. The labels C and D are as for ResNet; in addition, A refers to using the accumulated variance instead of the one-step predicted variance x_{t-1} . Thus, A + C + D is our approach. As before, we omit results for the ablation using the VC generalization bound since β is so small that the bound does not hold for any of the given α and β . In (b), we show the same ablations over the entire trajectory until $t=20$. As can be seen, using the calibrated predictor produces a large gain; these gains are most noticeable in the tails. Using the accumulated confidence produces a smaller, but still significant, gain. In (c) and (d), we show how the sizes vary with α and β , respectively. The trends are similar those for ResNet.

5 CONCLUSION

We have proposed an algorithm for constructing PAC confidence sets for deep neural networks. Our approach leverages statistical learning theory to obtain theoretical guarantees on the predicted confidence sets. These confidence sets quantify the uncertainty of deep neural networks. For instance, they can be used to inform safety-critical decision-making, and to ensure safety with high-probability in robotics control settings that leverage deep neural networks for perception. Future work includes extending these results to more complex tasks (e.g., structured prediction), and handling covariate shift (e.g., to handle policy updates in reinforcement learning).

ACKNOWLEDGMENTS

This work was supported in part by NSF CCF-1910769 and by the Air Force Research Laboratory and the Defense Advanced Research Projects Agency under Contract No. FA8750-18-C-0090. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Air Force Research Laboratory (AFRL), the Defense Advanced Research Projects Agency (DARPA), the Department of Defense, or the United States Government.

REFERENCES

- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Predictive inference with the jackknife. *arXiv preprint arXiv:1905.02928*, 2019.
- Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based reinforcement learning with stability guarantees. *Advances in neural information processing systems*, pp. 908–918, 2017.
- Marko Bohanec and Vladislav Rajkovic. Knowledge acquisition and explanation for multi-attribute decision making. *18th Intl Workshop on Expert Systems and their Applications*, 1998.
- Christopher P Bona de, A Russell Localio, John H Holmes, Vinay M Nadkarni, Shannon Stemler, Matthew MacMurchy, Miriam Zander, Kathryn E Roberts, Richard Lin, and Ron Keren. Video analysis of factors associated with response time to physiologic monitor alarms in a childrens hospital. *JAMA pediatrics* 171(6):524–531, 2017.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in Neural Information Processing Systems*, pp. 4754–4765, 2018.
- Paulo Cortez and Alice Maria Goncalves Silva. Using data mining to predict secondary school student performance. 2008.
- Morris H DeGroot and Stephen E Fienberg. The comparison and evaluation of forecasts of the Royal Statistical Society: Series D (The Statistician) 1(1-2):12–22, 1983.
- Christophe Denis and Mohamed Hebiri. Confidence sets with expected sizes for multiclass classification. *The Journal of Machine Learning Research* 18(1):3571–3598, 2017.
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1183–1192. JMLR. org, 2017.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. *arXiv preprint arXiv:1706.04592*, 2017.
- H Altay Guvenir, Burak Acar, Gulsen Demiroz, and Ayhan Cekin. A supervised machine learning algorithm for arrhythmia analysis. *Computers in Cardiology* 1997, pp. 433–436. IEEE, 1997.
- 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 (CVPR)*, pp. 770–778, 2016.
- David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 fps with deep regression networks. In *European Conference on Computer Vision*, pp. 749–765. Springer, 2016.
- Michael J Kearns and Umesh Virkumar Vazirani. *An introduction to computational learning theory*. MIT press, 1994.
- Alex Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv preprint arXiv:1404.5997*, 2014.
- Volodymyr Kuleshov and Percy S Liang. Calibrated structured prediction. *Advances in Neural Information Processing Systems*, pp. 3474–3482, 2015.
- Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. *arXiv preprint arXiv:1807.00263*, 2018.
- Jing Lei. Classification with confidence. *Biometrika* 101(4):755–769, 2014.

- Jing Lei, Max G. Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association* 113(523):1094–1111, 2018.
- Ali Malik, Volodymyr Kuleshov, Jiaming Song, Danny Nemer, Harlan Seymour, and Stefano Ermon. Calibrated model-based deep reinforcement learning. *International Conference on Machine Learning* pp. 4314–4323, 2019.
- Allan H. Murphy. Scalar and vector partitions of the probability score: Part i. two-state situation. *Journal of Applied Meteorology* 11(2):273–282, 1972.
- Mahdi Pakdaman Naeni, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. *Twenty-Ninth AAAI Conference on Artificial Intelligence* 2015.
- Harris Papadopoulos. Inductive conformal prediction: Theory and application to neural networks. In *Tools in artificial intelligence* IntechOpen, 2008.
- T. Pearce, M. Zaki, A. Brintrup, and A. Neely. High-quality prediction intervals for deep learning: A distribution-free, ensembled approach. *35th International Conference on Machine Learning, ICML 2018* volume 9, pp. 6473–6482, 2018.
- John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers* 1999.
- J. Ross Quinlan. Combining instance-based and model-based learning. *Proceedings of the Tenth International Conference on Machine Learning* pp. 236–243. Morgan Kaufmann Publishers Inc., 1993.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision* 115(3):211–252, 2015.
- Lukas Steinberger and Hannes Leeb. Leave-one-out prediction intervals in linear regression models with many variables. *arXiv preprint arXiv:1602.05801*, 2016.
- Lukas Steinberger and Hannes Leeb. Conditional predictive inference for high-dimensional stable algorithms. *arXiv preprint arXiv:1809.01412*, 2018.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *International Conference on Learning Representations (ICLR)* 2014.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* pp. 1–9, 2015.
- Leslie G. Valiant. A theory of the learnable. *Communications of the ACM* 27(11):1134–1142, 1984.
- Vladimir N. Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks* 10(5):988–999, 1999.
- Vladimir Vovk. Conditional validity of inductive conformal predictors. *Machine learning* 92(2-3): 349–376, 2013.
- Wenbo Wang and Xingye Qiao. Learning confidence sets using support vector machines. In *Advances in Neural Information Processing Systems* pp. 4929–4938, 2018.
- Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Online object tracking: A benchmark. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 2013.
- Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *Proceedings of the Eighteenth International Conference on Machine Learning* Citeseer, 2001.

Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 694–699. ACM, 2002.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. ICLR, 2017.

A DISCUSSION OF ALGORITHM DESIGN CHOICES

A.1 USEFULNESS OF TEMPERATURE SCALING

In this section, we discuss why temperature scaling can help improve the predicted confidence sets. A concern is that temperature scaling does not change the ordering of label probabilities. Thus, we may expect that temperature scaling does not affect the predicted confidence sets. However, this fact only holds when considering a single input—i.e., the ordering of the probabilities $f(y|x)$ for $y \in Y$ is not changed by temperature scaling. Indeed, the order of confidence sets for labels for different inputs can change. For a concrete example, consider two inputs x^0 , and the case $Y = \{0, 1, 2\}$. Assume that the label probabilities are

$$f(y|x) = [1/3 \quad 1/3 \quad 1/3]^T$$

$$f(y|x^0) = [3/4 \quad 1/4 \quad 0]^T$$

Now, if we take temperature very large, then the labels become roughly

$$f(y|x) = [1/3 \quad 1/3 \quad 1/3]^T$$

$$f(y|x^0) = [1/2 \quad 1/2 \quad 0]^T$$

As a consequence, there are confidence sets that are achievable when using f but are not achievable when using f^T . In particular, the confidence sets

$$C_T(x) = \{0\}$$

$$C_T(x^0) = \{0, 1\}$$

can be achieved using f (e.g., with $T = 2/5$). However, it is impossible to achieve these confidence sets using f^T for any choice of T , since if $1 \in C_T(x^0)$, then it must be the case that $C_T(x) = \{0, 1, 2\}$. Intuitively, we expect calibrated prediction to improve the ordering of probabilities across different inputs. Our experiments support this intuition, since they show that empirically, using calibrated predictors produces confidence sets of smaller size.

A.2 USEFULNESS OF DIRECT BOUND

One key design choice is to use a specialized generalization bound that directly provides PAC guarantees on our confidence sets rather than simply applying the VC dimension bound. The easiest way to determine which bound is better is to examine which one produces a smaller confidence set. In our approach, the size of the confidence set decreases monotonically with the choice of $(n; \epsilon)$ in (3). Thus, the bound that produces larger n is better. Recall that the VC dimension bound produces

$$n_{VC}(n; \epsilon) = \frac{\log(2n) + 1}{\epsilon} \frac{\log(4/\epsilon)}{\epsilon},$$

whereas our direct bound produces (for $\epsilon > 0$)

$$n_{\text{direct}}(n; \epsilon) = \max_{k \geq 2N} k \text{ subj. to } \sum_{i=0}^k \binom{n}{i} (1-\epsilon)^n < \epsilon$$

Directly comparing these two choices of n is difficult, but our experiments show empirically that using the direct bound outperforms using the indirect bound.

A more direct way to compare the two approaches is to instead ask how large n needs to be to achieve $(n; \epsilon) = 0$. For n_{VC} , it is easy to check that we need

$$n \geq \frac{\log(2n) + 1 + \log(4/\epsilon)}{\epsilon}.$$

Thus, we need n to be at least $O(\log(1/\epsilon)^2)$ (and possibly greater, to account for the $\log(2n)$ term). In contrast, for our direct bound, $n = 0$ corresponds to the case $\epsilon = 0$. To achieve $n = 0$, it suffices to have n satisfying $(1-\epsilon)^n < \epsilon$. Using $(1-\epsilon)^n \leq e^{-n\epsilon}$, it suffices to have n satisfying

$$n \geq \frac{\log(1/\epsilon)}{\epsilon}.$$

Figure 4: Sample complexity of different bounds; we $\epsilon = 10^{-5}$. Left: Sample complexity of VC bound and direct bound when $\epsilon = 0$. Right: Sample complexity of direct bound for varying

In other words, n only needs to be $\mathcal{O}(\log(1/\epsilon))$. For small ϵ (e.g., $\epsilon = 0.01$), we need 100 fewer samples to achieve the same size confidence set (i.e., with $\epsilon = 0$). In Figure 4 (right), we compute the exact values needed to get $(n; \epsilon) = 0$ as a function of ϵ for each bound ($\epsilon = 10^{-5}$). As expected, our bound requires substantially smaller

Finally, in Figure 4 (right), we compare the magnitude needed to achieve larger values of ϵ using our direct bound; for simplicity, we actually consider larger values of ϵ (where $\epsilon = k/n$), but the qualitative insights are the same. As can be seen, even for large k (e.g., $k = 50$), the number of samples increases, but not substantially.

B THEORETICAL GUARANTEES

B.1 ASSUMPTIONS

We make two additional technical assumptions in Theorem 1, both of which are standard. First, we assume that f is measurable; this assumption holds for all models used in practice, including neural networks (e.g., it holds as long as f is continuous).

Second, letting $\mu : Z \rightarrow \mathbb{R}$, where $Z = X \times Y$, be defined by $\mu((x; y)) = -\log f(y | x)$, we assume that the distribution \mathbb{D} induced by μ on R has continuous cumulative distribution function (CDF). More precisely, letting \mathbb{D} be the measure defined on R , then \mathbb{D} is defined by the measure

$$\mathbb{D}(t) = \mathbb{D}(\mu^{-1}(t));$$

where $\mu^{-1} : \mathbb{R} \rightarrow 2^Z$ is the inverse of μ in the sense that $z \in \mu^{-1}(\mu(z))$ for all $z \in Z$. Then, we assume that the CDF corresponding to \mathbb{D} is continuous. This second assumption is standard in statistical learning theory (Kearns & Vazirani, 1994). Essentially, it says that for any R , the probability that $-\log f(y | x)$ must equal zero. This assumption should hold unless $f(y | x)$ or $f(y | x)$ are degenerate in some way. Furthermore, we can detect this case. In particular, the failure mode corresponds to the case that we see multiple points with the same value $-\log f(y | x)$. Thus, choosing $\epsilon = -\log f(y | x)$ would include all these points, so the realized error rate is larger than desired for ϵ . In this case, we can simply choose a slightly larger ϵ to avoid this problem.

B.2 PROOF OF THEOREM 1

At a high level, our proof proceeds in three steps. First, we show that a confidence set predictor can be encoded as a binary classifier M_T . Second, we show that a PAC bound M_T implies a PAC bound for C_T (where in both cases, the unknown parameter $\theta \in R$). Third, we prove PAC bounds on the error M_T ; by the second step, these bounds complete our proof.

Encoding C_T as a binary classifier M_T . We begin by showing how the problem of learning a PAC confidence set predictor C_T reduces to the problem of learning a PAC binary classifier M_T . First, we show that for any $\theta \in R$, the confidence set predictor C_T can be encoded as a binary classifier M_T . Consider any parameter $\theta \in R$. Recall that we use the model $f(y | x)$ to construct the confidence set predictor

$$C_T(x) = \{y \in Y \mid f(y | x) \geq \epsilon\};$$

Now, define the map $\psi : Z \rightarrow \mathbb{R}$ by $\psi(x; y) = -\log f(y | x)$, where $Z = X \times Y$, and define the binary classifier $M_T : \mathbb{R} \rightarrow \{0, 1\}$ by

$$M_T(t) = \mathbb{I}[t \leq T].$$

Here, $\mathbb{I}[s]$ is the indicator function, which returns one if a statement is true and zero otherwise. We claim that

$$C_T(x) = \int_Y \int_X M_T(\psi(x; y)) f(y | x) dP(y | x) = 1 - g(x) \quad (12)$$

To see this claim, note that

$$\begin{aligned} C_T(x) &= \int_Y \int_X M_T(\psi(x; y)) f(y | x) dP(y | x) \\ &= \int_Y \int_X \mathbb{I}[\psi(x; y) \leq T] f(y | x) dP(y | x) \\ &= \int_Y \int_X \mathbb{I}[\log f(y | x) \geq -T] f(y | x) dP(y | x) \\ &= \int_Y \int_X \mathbb{I}[f(y | x) \geq e^{-T}] f(y | x) dP(y | x) \\ &= \int_Y \int_X M_T(\psi(x; y)) f(y | x) dP(y | x) = 1 - g(x) \end{aligned}$$

as claimed.

PAC bound for M_T implies PAC bound for C_T . Next, we show that a PAC bound for M_T implies a PAC bound for C_T . More precisely, we design a data distribution \mathcal{D} and loss ℓ , and show that (i) the distribution of \hat{F} (trained to optimize M_T) is the same as the distribution \hat{C} (constructed using our algorithm), and (ii) a PAC bound for M_T (where \hat{F} is trained on data from \mathcal{D}) implies a PAC bound for C_T . We show that as a consequence, a PAC bound for M_T implies a PAC bound for C_T .

We begin by constructing \mathcal{D} and ℓ . To this end, recall that D is a given distribution over $X \times Y$. We define a data distribution \mathcal{D} over $X \times Y$, where $X = \mathbb{R}$ and $Y = \{0, 1\}$, as follows. The first component of \mathcal{D} is the distribution over X induced by D , and the second component is the distribution over Y that places all probability mass on $\{0, 1\}$. Formally, \mathcal{D} exists assuming ψ is measurable, so the induced distribution exists; for all our choices of ψ (i.e., categorical or Gaussian), this property is satisfied. Then,

$$\mathcal{D}((t; a)) = D(\psi^{-1}(t)) \mathbb{I}[a = 1];$$

where D is the measure encoding D , and \mathbb{I} is the measure encoding \mathbb{I} . Furthermore, we define $\ell : Y \times Y \rightarrow \{0, 1\}$ to be the 0-1 loss $\ell(a; a^0) = \mathbb{I}[a \neq a^0]$. Finally, let \hat{F} be chosen using our algorithm—i.e.,

$$\begin{aligned} \hat{F} &= \arg \min_T T \text{ subj. to } L(C_T; Z) \\ L(C_T; Z) &= \frac{1}{|Z|} \sum_{(x; y) \in Z} \mathbb{I}[y \neq C_T(x)]; \end{aligned}$$

for any $\epsilon \in \mathbb{R}_{>0}$, and let \hat{F} be chosen similarly for M_T —i.e.,

$$\begin{aligned} \hat{F} &= \arg \min_T T \text{ subj. to } L(M_T; Z) \\ L(M_T; Z) &= \frac{1}{|Z|} \sum_{(t; a) \in Z} \mathbb{I}[M_T(t) \neq a]; \end{aligned}$$

Now, we show (i) above. In particular, we claim that $\hat{F}(Z)$ has the same distribution as $\hat{C}(Z)$, where $Z \subset \mathcal{D}^n$ and $Z \subset \mathcal{D}^n$ are random datasets. To this end, define $Z \subset \mathcal{D}^n$ by

$$Z = ((z_1; 1), \dots, (z_n; 1));$$

Note that

$$\begin{aligned} L(M_T; Z) &= \frac{1}{|Z|} \sum_{i=1}^n \mathbb{I}[M_T(\psi(x_i; y_i)) \neq 1] \\ &= \frac{1}{|Z|} \sum_{i=1}^n \mathbb{I}[y_i \neq C_T(x_i)] \\ &= L(C_T; Z); \end{aligned}$$

from which it follows that

$$\begin{aligned} \hat{T}(Z) &= \arg \min_T \text{ subj. to } L(C_T; Z) \\ &= \arg \min_T \text{ subj. to } \mathbb{E}(M_T; (Z)) \\ &= T((Z)): \end{aligned}$$

By construction of \hat{T} , the random variables $\hat{T}(Z)$ and $T((Z))$ have the same distribution; thus, it follows that the random variables $\mathbb{E}(\hat{T}(Z))$ and $\mathbb{E}(T((Z)))$ have the same distribution as well. Since $\mathbb{E}(\hat{T}(Z)) = \mathbb{E}(T((Z)))$, it follows that $\hat{T}(Z)$ has the same distribution as $T((Z))$, as claimed.

Next, we show (ii) above. In particular, we claim that a PAC bound for $\mathbb{E}(M_{T((Z))})$ —i.e.,

$$P_{Z \sim \mathcal{D}^n} [\mathbb{E}(M_{T((Z))}) \leq \epsilon] \geq 1 - \delta;$$

implies a PAC bound for $\mathbb{E}(L(C_{\hat{T}(Z)}))$ —i.e.,

$$P_{Z \sim \mathcal{D}^n} [L(C_{\hat{T}(Z)}) \leq \epsilon] \geq 1 - \delta;$$

where the true losses are

$$\begin{aligned} \mathbb{E}(M_T) &= \mathbb{E}_{(t;a) \sim \mathcal{D}} [M_T(t; a)] = P_{(t;a) \sim \mathcal{D}} [M_T(t) \leq a] \\ L(C_T) &= \mathbb{E}_{(x;y) \sim \mathcal{D}} [I[y \geq C_T(x)]] = P_{(x;y) \sim \mathcal{D}} [y \geq C_T(x)]; \end{aligned}$$

Note that it suffices to show that the true loss for \hat{T} equals the true loss for T —i.e.,

$$L(C_{\hat{T}}) = \mathbb{E}(M_T);$$

since this equation (together with the PAC bound for $\mathbb{E}(M_{T((Z))})$) implies

$$P_{Z \sim \mathcal{D}^n} [L(C_{\hat{T}(Z)}) \leq \epsilon] = P_{Z \sim \mathcal{D}^n} [\mathbb{E}(M_{T((Z))}) \leq \epsilon] \geq 1 - \delta;$$

as desired. To see the claim, note that

$$\begin{aligned} \mathbb{E}(M_T) &= \int_{\mathcal{Z}} P_{(t;a) \sim \mathcal{D}} [M_T(t) \leq a] \\ &= \int_{\mathcal{Z}} I[M_T(t) \leq a] d_{\mathcal{D}}((t; a)) \\ &= \int_{\mathcal{Z}} I[a = 1] \int_{\mathcal{X}} I[M_T(t) \leq a] d_{\mathcal{D}}(t) \\ &= \int_{\mathcal{Z}} I[M_T(t) \leq 1] d_{\mathcal{D}}(t) \end{aligned}$$

Now, using the change of variables $(z) = (x; y)$, we have

$$\begin{aligned} \mathbb{E}(M_T) &= \int_{\mathcal{Z}} I[M_T((z)) \leq 1] d_{\mathcal{D}}(z) \\ &= \int I[M_T((x; y)) \leq 1] D(x; y) dx dy; \end{aligned}$$

Then, using (12), we have

$$\begin{aligned} \mathbb{E}(M_T) &= \int I[y \geq C_T(x)] D(x; y) dx dy \\ &= P_{(x;y) \sim \mathcal{D}} [y \geq C_T(x)] \\ &= L(C_T); \end{aligned}$$

as claimed.

Finally, combining (i) and (ii), we have

$$P_{Z \sim \mathcal{D}^n} [L(C_{\hat{T}(Z)}) \leq \epsilon] = P_{Z \sim \mathcal{D}^n} [L(C_{T((Z))}) \leq \epsilon] \geq 1 - \delta;$$

where the first equality follows since (i) says that $\hat{A}(Z)$ (where $Z \in \mathcal{D}^n$) has the same distribution as $\mathcal{T}(Z)$ (where $Z \in \mathcal{D}^n$), and the second inequality follows by (ii).

Generalization bound. Finally, we prove the PAC bound

$$P_{Z \in \mathcal{D}^n} [\mathcal{L}(M_{\mathcal{T}}) \leq \epsilon] \geq 1 - \epsilon; \quad (13)$$

for $M_{\mathcal{T}}$, where $\epsilon = \sum_{i=0}^k \binom{n}{i} (1 - \epsilon)^n$; for conciseness, we have dropped the dependence of \mathcal{T} on Z . By the previous step, this bound implies the theorem statement. To this end, we first simplify the left-hand side of the inequality (13). In particular, let \mathcal{T} be the smallest \mathcal{T} for which $\mathcal{L}(M_{\mathcal{T}}) = \epsilon$; such a \mathcal{T} exists by our assumption that \mathcal{P} has continuous density function.

First, we claim that $\mathcal{T} < T$ implies $\mathcal{L}(M_{\mathcal{T}}) > \mathcal{L}(M_T)$. Assuming $\mathcal{T} < T$, then

$$\begin{aligned} \mathcal{L}(M_{\mathcal{T}}) &= P_{(t;a) \in \mathcal{D}} [M_{\mathcal{T}}(t) \leq a] \\ &= E_{(t;a) \in \mathcal{D}} [I[M_{\mathcal{T}}(t) \leq a]] \\ &= E_{(t;a) \in \mathcal{D}} [I[M_{\mathcal{T}}(t) \leq 1]] \\ &= E_{(t;a) \in \mathcal{D}} [I[t \leq T]] \\ &= E_{(t;a) \in \mathcal{D}} [I[t > T]] \\ &> E_{(t;a) \in \mathcal{D}} [I[t > T]] \\ &= \mathcal{L}(M_T); \end{aligned}$$

Assuming $\mathcal{T} \geq T$, we can similarly show that $\mathcal{L}(M_{\mathcal{T}}) \leq \mathcal{L}(M_T)$. It follows that

$$\begin{aligned} P_{Z \in \mathcal{D}^n} [\mathcal{L}(M_{\mathcal{T}}) > \epsilon] &= P_{Z \in \mathcal{D}^n} [\mathcal{L}(M_{\mathcal{T}}) > \mathcal{L}(M_T)] \\ &= P_{Z \in \mathcal{D}^n} [\mathcal{T} < T]; \end{aligned}$$

As a consequence, (13) is equivalent to

$$P_{Z \in \mathcal{D}^n} [\mathcal{T} < T] \leq \epsilon;$$

Next, recall that \mathcal{T} must satisfy $\mathcal{L}(M_{\mathcal{T}}; Z) = \epsilon$, where

$$\mathcal{L}(M_{\mathcal{T}}; Z) = \frac{1}{n} \sum_{(t;a) \in \mathcal{D}} I[M_{\mathcal{T}}(t) \leq a];$$

Assuming $\mathcal{T} < T$, and using $k = n$, it follows that

$$\begin{aligned} \sum_{(t;a) \in \mathcal{D}} I[M_{\mathcal{T}}(t) \leq a] &= \sum_{(t;a) \in \mathcal{D}} I[M_{\mathcal{T}}(t) \leq 1] \\ &= \sum_{(t;a) \in \mathcal{D}} I[t > T] \\ &= \sum_{(t;a) \in \mathcal{D}} I[t > T]; \end{aligned}$$

As a consequence, we have

$$\begin{aligned} P_{Z \in \mathcal{D}^n} [\mathcal{T} < T] &= P_{Z \in \mathcal{D}^n} \left[\sum_{(t;a) \in \mathcal{D}} I[t > T] \geq k \right] \\ &= \sum_{i=0}^k P_{Z \in \mathcal{D}^n} \left[\sum_{(t;a) \in \mathcal{D}} I[t > T] \geq i \right] = \sum_{i=0}^k \epsilon^i; \end{aligned}$$

By our definition of T , the event in the final expression says that the sum of n i.i.d. Bernoulli random variables $[t > T]$ Bernoulli(p) is at most k . Thus, this event follows a distribution $\text{Binomial}(n; p)$, so

$$P_{Z \sim \mathcal{D}^n} [T < T] = \sum_{i=0}^k \binom{n}{i} p^i (1-p)^{n-i} = \epsilon;$$

as claimed. The theorem statement follows.

C DETAILS ON PROBABILITY FORECASTERS FOR SPECIFIC TASKS

In this section, we describe architectures for probability forecasters for classification, regression, and model-based reinforcement learning.

Classification. For the case $\mathcal{Y} = \{1, \dots, Y\}$, we choose the probability forecast to be a neural network with a softmax output. Then, we can compute a confidence set

$$C_T(x) = \{y \in \mathcal{Y} \mid f(y|x) \geq \epsilon\}$$

by explicitly enumerating \mathcal{Y} . We measure the size $\mathcal{O}_T(x)$ as $|C_T(x)|$.

Regression. For the case $\mathcal{Y} = \mathbb{R}$, we choose the probability forecast to be a neural network that outputs the parameters (μ, Σ) of a Gaussian distribution. Then, we have

$$C_T(x) = \left\{ \frac{\mu}{\sqrt{2(T - \log(\frac{p}{2}))}}; \pm \frac{\Sigma}{\sqrt{2(T - \log(\frac{p}{2}))}} \right\}$$

This choice generalizes $\mathcal{C}_T = \mathbb{R}^d$ by having f output the parameters (μ, Σ) of a d -dimensional Gaussian distribution. Note that $C_T(x)$ is an ellipsoid $C_T(x) = \mu + \Sigma^{1/2} S^d$, where $\mu \in \mathbb{R}^d$ and S^d is the unit sphere in \mathbb{R}^d ; in particular, $\Sigma = D^{-1/2} Q$, where $Q D Q^T$ is the eigendecomposition of Σ .

$$(2T - d \ln 2 - \ln \det \Sigma)^{-1/2}$$

We measure the size $\mathcal{O}_T(x)$ as $\|k_F\|_F$, where $\|k_F\|_F$ is the Frobenius norm.

Model-based reinforcement learning. In model-based reinforcement learning, the goal is to predict trajectories based on a model of the dynamics. We consider an MDP with state $\mathcal{X} \subseteq \mathbb{R}^{d_x}$, actions $U \subseteq \mathbb{R}^{d_u}$, an unknown distribution over initial states \mathcal{S}_0 , and unknown dynamics $(x^0_j | x; u)$ mapping a state-action pair $(x; u) \in \mathcal{X} \times U$ to a distribution over states \mathcal{X} . We assume a fixed, known policy $(u | x)$, mapping a state $x \in \mathcal{X}$ to a distribution over actions U . The (unknown) closed-loop dynamics are $f(x^0_j | x) = E_{(u|x)}[g(x^0_j | x; u)]$.

Given initial state $x_0 \in \mathcal{X}$ and time horizon $H \in \mathbb{N}$, we can sample a trajectory $x_{1:H} = (x_1; \dots; x_H)$ by setting $x_0 = x_0$ and sequentially sampling $x_{t+1} \sim f(\cdot | x_t)$ for $t \in \{0, 1, \dots, H-1\}$. Our goal is to predict a confidence set $\mathcal{C}_T(x_0) \subseteq \mathcal{X}^H$ that contains $x_{1:H} \in \mathcal{X}^H$ with high probability (according to both the randomness in initial states \mathcal{S}_0 and f). This problem is a multivariate regression problem with input x_0 and output $\mathcal{Y} = \mathcal{X}^H$.

We assume given a probability forecaster $f(x^0_j | x) = N(x^0_j; \mu(x); \Sigma(x))$ trained to predict the distribution over next states—i.e., $f(x^0_j | x) = f(x^0_j | x)$. Given initial state $x_0 \in \mathcal{X}$ and time horizon $H \in \mathbb{N}$, we construct the mean trajectory $\bar{x}_{1:H}$ by setting $x_0 = x_0$ and letting $x_{t+1} = \bar{x}_t$. To account for the fact that the variances accumulate over time, we sum them together to obtain the predicted variances $\tilde{\Sigma}_{1:H}$ —i.e.,

$$\tilde{\Sigma}_t = \Sigma(x_0) + \Sigma(x_1) + \dots + \Sigma(x_{t-1});$$

Then, we use the probability forecaster $f(x_{1:H}; \tilde{\Sigma}_{1:H}) = N(x_{1:H}; \bar{x}_{1:H}; \tilde{\Sigma}_{1:H})$ (where we think of $x_{1:H}$ as a vector in $\mathbb{R}^{H \cdot d_x}$ and $\tilde{\Sigma}_{1:H}$ as a block diagonal matrix in $\mathbb{R}^{(H \cdot d_x) \times (H \cdot d_x)}$) to construct confidence sets.

Finally, we describe how we measure the size of a predicted confidence set $\mathcal{C}_T(x_0) \subseteq \mathcal{X}^H$. In particular, note that $\mathcal{C}_T(x_0)$ has the form

$$C_T(x_0) = (C_{T;1}(x_0); \dots; C_{T;H}(x_0));$$

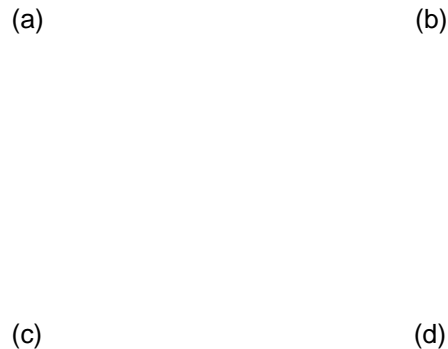


Figure 5: Comparison to baselines that do not have theoretical guarantees. In (a) and (b), we show results for ImageNet, and in (c) and (d), we show results for the half-cheetah. In (a) and (c), we show the empirical error in the confidence set sizes; the dotted line denotes 0.01, our target confidence set error. In (b) and (d), we show the sizes of the constructed confidence sets.

i.e., $C_{T,t}(x_0)$ is the confidence set for the state reached after time steps. Then, we measure the size of the confidence set for each component $C_{t,k}(x_0)$ (for $t \in \{1, \dots, H\}$) individually, and take the average. As in the case of regression, $C_{T,t}(x_0)$ is an ellipsoid $C_{T,t}(x_0) = x_t + \sum_{k=1}^d \sigma_{t,k}^2 e_k e_k^T$; then, the size of $C_T(x_0)$ is $\sum_{t=1}^H \sum_{k=1}^d \sigma_{t,k}^2$.

An additional detail is that when we calibrate this forecaster, we calibrate each component individually—i.e., we use dH calibration parameters; $\dots; H$.

D ADDITIONAL RESULTS

D.1 COMPARISON TO ADDITIONAL BASELINES

We compare to two baselines that do not have theoretical guarantees. We assume given a probability forecast $f(y_j|x)$. Then, given an input $x \in \mathcal{X}$, we construct the confidence set to satisfy

$$\sum_{y \in C(x)} f(y_j|x) \geq 1 - \epsilon \quad (14)$$

More precisely, we first rank the labels in decreasing order of $f(y_j|x)$, to obtain a list $(y_1; y_2; \dots; y_{|Y|})$. Then, we choose the smallest k such that (14) holds for $C(x) = \{y_1; \dots; y_k\}$. Intuitively, if the probabilities $f(y_j|x)$ are correct (i.e. $f(y_j|x)$ is the true probability of y_j given x), then this confidence set should contain the true label with high probability.

For regression, we cannot explicitly rank labels $Y \subseteq \mathbb{R}^d$, but they are monotonically decreasing away from the mean. Then, assuming $f(y_j|x) = \mathcal{N}(y_j; \mu(x); \Sigma(x))$ is Gaussian, we take an ellipsoid of shape $C(x)$ around $\mu(x)$ with minimum radius that captures $1 - \epsilon$ of the probability mass of $f(y_j|x)$. More precisely, we choose

$$\begin{aligned} C(x) &= C_{T(x)}(x) \\ \hat{T}(x) &= \arg \min_{T \subseteq \mathbb{R}^d} T \text{ subj. to } P_{f(y_j|x)}[y \in C_T(x)] \geq 1 - \epsilon; \end{aligned}$$

Figure 6: Confidence set sizes for two neural network architectures trained on ImageNet; for both, we use $n = 20,000$, $\epsilon = 0.01$ and $\delta = 10^{-5}$. Left: AlexNet (Krizhevsky, 2014); here, the empirical confidence set error of our approach is 0.0066. Right: GoogLeNet (Szegedy et al., 2015); here, the empirical confidence set error of our approach is 0.0061.

where $C_T(x) = \{y \in Y \mid f(y|x) \geq \tau\}$ as before. Note that unlike our algorithm, the threshold $\hat{\tau}(x)$ is not a learned parameter, but is computed independently for each new input x . We can solve for $\hat{\tau}(x)$ efficiently by changing basis to convert $f(\cdot|x)$ to a standard Gaussian distribution, and then using the error function to compute the cutoff that includes the desired probability mass.

In Figure 5, we compare the confidence sets constructed using this approach with (i) the forecaster $f_\lambda(y|x)$ without any calibration, and (ii) the calibrated forecaster $f_{\lambda^*}(y|x)$. We plot both the confidence set sizes and the empirical error rates. For the latter, recall that a confidence set predictor C is correct if $L(C) < \epsilon$, where $L(C)$ the true error rate. However, we cannot measure $L(C)$; instead, we approximate it on a held-out test set $Z_{\text{test}} \sim X \times Y$ —i.e., $L(C) \approx \hat{L}(C; Z_{\text{test}})$, where

$$\hat{L}(C; Z_{\text{test}}) = \frac{1}{|Z_{\text{test}}|} \sum_{(x,y) \in Z_{\text{test}}} \mathbb{1}[y \notin C(x)].$$

Intuitively, $\hat{L}(C; Z_{\text{test}})$ is the fraction of inputs $(x,y) \in Z_{\text{test}}$ such that the predicted confidence set for x does not contain y . We say a confidence set is empirically valid when $\hat{L}(C; Z_{\text{test}}) < \epsilon$. Recall that our algorithm guarantees correctness with probability at least $1 - \delta$ where $\delta = 10^{-5}$.

As can be seen, the baseline approaches are not empirically valid in all cases. In one case—namely, the baseline with the calibrated forecaster on ImageNet—the confidence sets are almost empirically valid. However, in this case, the confidence sets are much larger than those based on our approach, despite the fact that the error rate of our confidence sets are empirically valid. Thus, our algorithms outperform the baselines in all cases.

D.2 RESULTS ON ADDITIONAL IMAGENET NEURAL NETWORK ARCHITECTURES

We apply our approach to two additional neural network architectures for ImageNet: AlexNet (Krizhevsky, 2014) and GoogLeNet (Szegedy et al., 2015). Our results are shown in Figure 6. As can be seen, calibration reduces the confidence set sizes for AlexNet, but actually increases the confidence set sizes for GoogleNet. Thus, both calibrated and uncalibrated models may need to be considered when constructing confidence set predictors. Also, we find that confidence set sizes are correlated with classification error—the test errors for AlexNet, GoogleNet, and ResNet-101 are 29.41%, 21.34%, and 17.03% respectively, and their confidence set sizes decrease in the same order.

D.3 RESULTS ON ADDITIONAL CLASSIFICATION DATASETS

We apply our approach to three small classification datasets: an Arrhythmia detection dataset (Guvénir et al., 1997), a car evaluation dataset (Bohanec & Rajkovic, 1988), and a medical alarm dataset (Bona de et al., 2017). The confidence set sizes are shown in Figure 7. We choose larger values of ϵ and δ since we cannot obtain confidence sets that satisfy the PAC criterion with smaller ϵ when the number of validation examples is too small. For all three datasets, the empirical confidence set error is smaller than the specified error ϵ , thus, the constructed confidence sets are empirically valid. For these datasets, the confidence set sizes of our approach and our approach without calibration are similar, most likely due to the small number of class labels.

(a) Arrhythmia

(b) Car

(c) CHOP Alarm

(d) Suppressed alarms

Figure 7: Confidence set sizes for three additional classification benchmarks: (a) the arrhythmia detection dataset (Guvenir et al., 1997); here $n = 90$, $\epsilon = 0.1$, $\delta = 0.05$, and the empirical confidence set error of our approach is 0.0435 (b) the car evaluation dataset (Bohanec & Rajkovic, 1988); here $n = 345$, $\epsilon = 0.05$, $\delta = 10^{-5}$, and the empirical confidence set error of our approach is 0.0172 and (c) the CHOP alarm dataset (Bona de et al., 2017); here $n = 1000$, $\epsilon = 0.02$, $\delta = 10^{-5}$, and the empirical confidence set error of our approach is 0.0159 (d) The fractions of actionable and false alarms with a confidence set (i.e., only contains false alarm).

We additionally ran our approach on a medical dataset where classification decisions are safety critical; thus, correct predicted confidence sets are required. In particular, we use the Children's Hospital of Philadelphia (CHOP) alarm dataset (Bona de et al., 2017). This dataset consists of vital signs from 100 patients around one year of age. One of the vital signs is the oxygen level of the blood, and a medical device generates an alarm if the oxygen level is below a specified level. The labels indicate whether the generated alarm is true ($y = 1$) or false ($y = 0$). We use $n = 1000$, $\epsilon = 0.02$, and $\delta = 10^{-5}$. The empirical confidence set error of our approach is $\hat{L}(C; Z_{\text{test}}) = 0.0159$.

The key question is how many false alarms can be reliably detected using machine learning to help reduce alarm fatigue. We consider an approach where we use the predicted confidence sets to detect false alarms. In particular, we first train a probability forecaster $X \rightarrow P_Y$, where $Y = \{0, 1\}$, to predict the probability that an alarm is true, and then construct a calibrated confidence set predictor $f: X \rightarrow \mathcal{Z}^Y$ based on this forecaster. We consider an alarm to be false if the predicted confidence set is $f(x) = \{0\}$ —i.e., according to our confidence set predictor, the alarm is definitely false. Then, our PAC guarantee says that the alarm is actually false with probability at least $1 - \epsilon$. In summary, we suppress an alarm if $f(x) = \{0\}$. Using our approach, 176=630 (i.e., 27.94%) of false alarms are suppressed, while only 187 (i.e., 6.95%) true alarms are suppressed (see Figure 7 (d)).

D.4 RESULTS ON ADDITIONAL REGRESSION DATASETS

We ran our algorithm on two small regression baselines—the Auto MPG dataset (Quinlan, 1993) and the student grade dataset (Cortez & Silva, 2008). We show results in Figure 8. The parameters we use are $\epsilon = 0.1$ and $\delta = 0.05$; as with the smaller classification datasets, we use larger choices of ϵ and δ since we cannot construct valid confidence sets for smaller choices. For the Auto MPG dataset, the empirical confidence set error of our model is $\hat{L}(C; Z_{\text{test}}) = 0.0597$, so these are empirically valid. For the student grade dataset, the empirical confidence set error is $\hat{L}(C; Z_{\text{test}}) = 0.1250$ which is slightly larger than desired; this failure is likely due to the fact that the failure probability $\epsilon = 0.05$ is somewhat large.

Figure 8: Confidence set sizes for two benchmarks focused on regression; for both, we use $\alpha = 0.05$. Left: the Auto MPG dataset (Quinlan, 1993); here, $n = 70$, and the empirical confidence set error of our approach is 0.1250. Right: The student grade dataset (Cortez & Silva, 2008); here, $n = 100$, and the empirical confidence set error of our approach is 0.0597.

D.5 ADDITIONAL RESULTS ON IMAGENET, HALF-CHEETAH, AND OBJECT TRACKING

Table 4 & 5 show examples of ResNet confidence set sizes for ImageNet images. Table 6 shows results for varying α on ResNet. Tables 7 & 8 show results for varying α on the Half-Cheetah. Table 9 shows visualizations of the confidence sets predicted for our object tracking benchmark.

	$ C(x) = 1$	$ C(x) = 5$	$ C(x) = 10$	$ C(x) = 50$	$ C(x) = 100$	$ C(x) = 200$
airship						
zebra						
banana						
carousel						
star sh						
street sign						

Table 4: ImageNet images with varying ResNet con dence set sizes. The con dence set sizes are on the top. The true label is on the left-hand side. Incorrectly labeled images are boxed in red.

$1 \leq j \leq 5$	$5 \leq j \leq 10$	$10 \leq j \leq 20$
king penguin	Chihuahua toy terrier Italian greyhound Boston bull miniature pinscher	banded gecko common iguana American chameleon whiptail; agama frilled lizard; alligator lizard green lizard African chameleon Komodo dragon
shopping basket	English springer Welsh springer spaniel collie; boxer Saint Bernard Leonberg	altar analog clock bell cote castle church cinema dome monastery palace vault wall clock
chambered nautilus	face powder hamper lotion; packet shopping basket	barber chair hand blower medicine chest paper towel plunger shower curtain soap dispenser toilet seat tub; washbasin washer toilet tissue
bonnet	kite; bald eagle vulture great grey owl bittern	beach wagon cab; car wheel convertible grille; limousine minivan mobile home passenger car pickup recreational vehicle sports cartow truck
Madagascar cat fadi	tiger cat lynx; leopard snow leopard jaguar tiger; cheetah	cannon castle cliff dwelling; megalith monastery obelisk prison stone wall triumphal arch vault alp; cliff; promontory valley
ballpoint fountain pen	cash machine desktop computer entertain. center home theater loudspeaker monitor; screen television	amphibian cassette player re engine minibus minivan passenger car pole; police van puck; racer radio school bus screwdriver streetcar trolleybus
ibex impala gazelle	common iguana whiptail; agama frilled lizard; alligator lizard green lizard Komodo dragon African crocodile American alligator	juncq water ouzel water snake drake red-breasted merganser goose cray sh; little blue heron European gallinule ruddy turnstone red-backed sandpiper redshank dowitcher oystercatcher albatross otter
indigo bunting bee eater hummingbird jacamar	barber chair barbershop electric fan hand blower iron; rocking chair table lamp tricycle; vacuum	accordion acoustic guitar banjo bassoon cornet drum drumstick electric guitar French horn maraca microphone oboe sax stage torch trumpet violin

Table 5: Condense sets of ImageNet images with varying ResNet condense set sizes. The predicted condense set is shown to the right of the corresponding input image. The true label is shown in red, and the predicted label is shown with a hat.

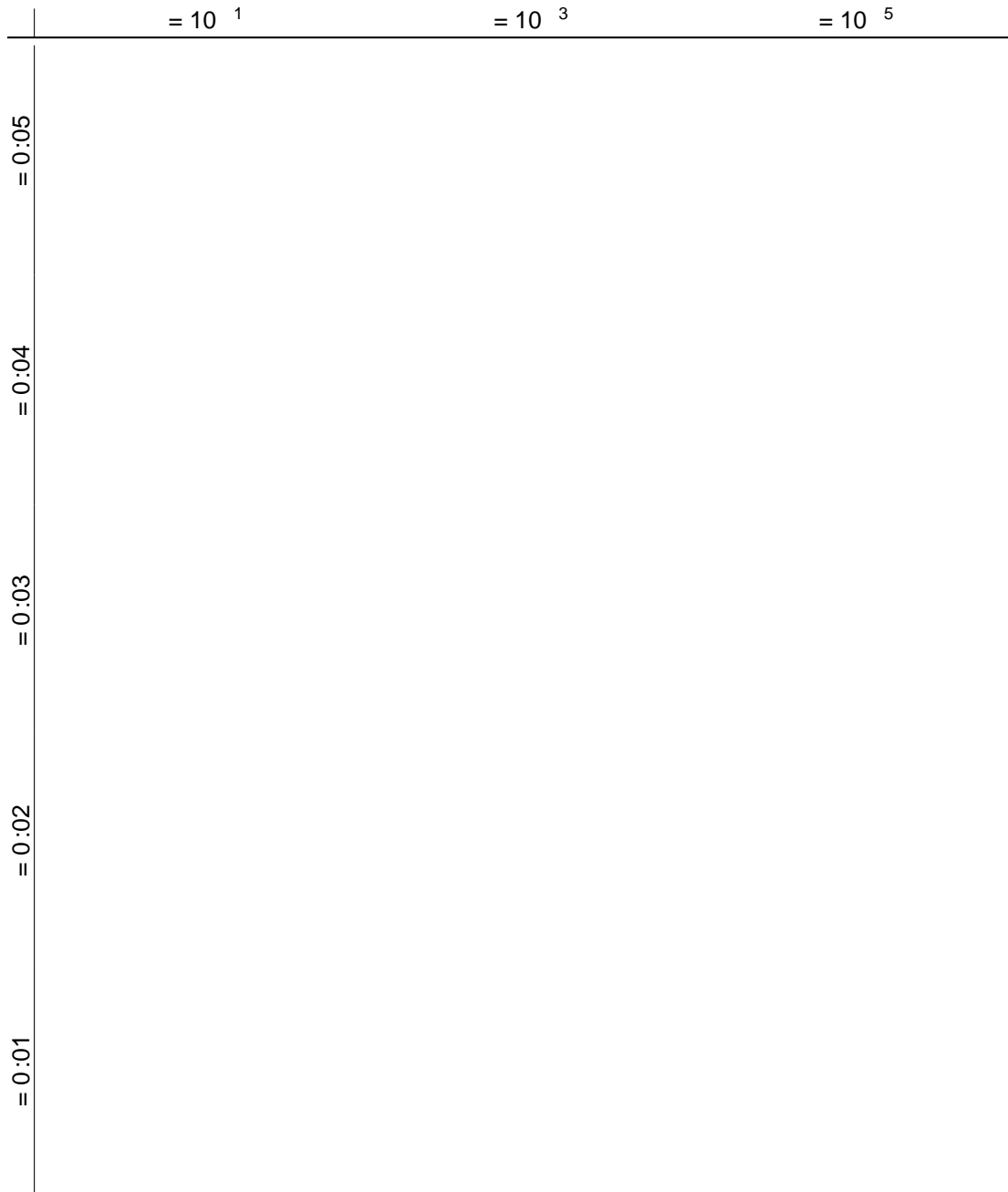


Table 6: Confidence set sizes for ResNet trained on ImageNet, for varying confidence level and sample size $n = 20,000$. The plots are as in Figure 1 (a).

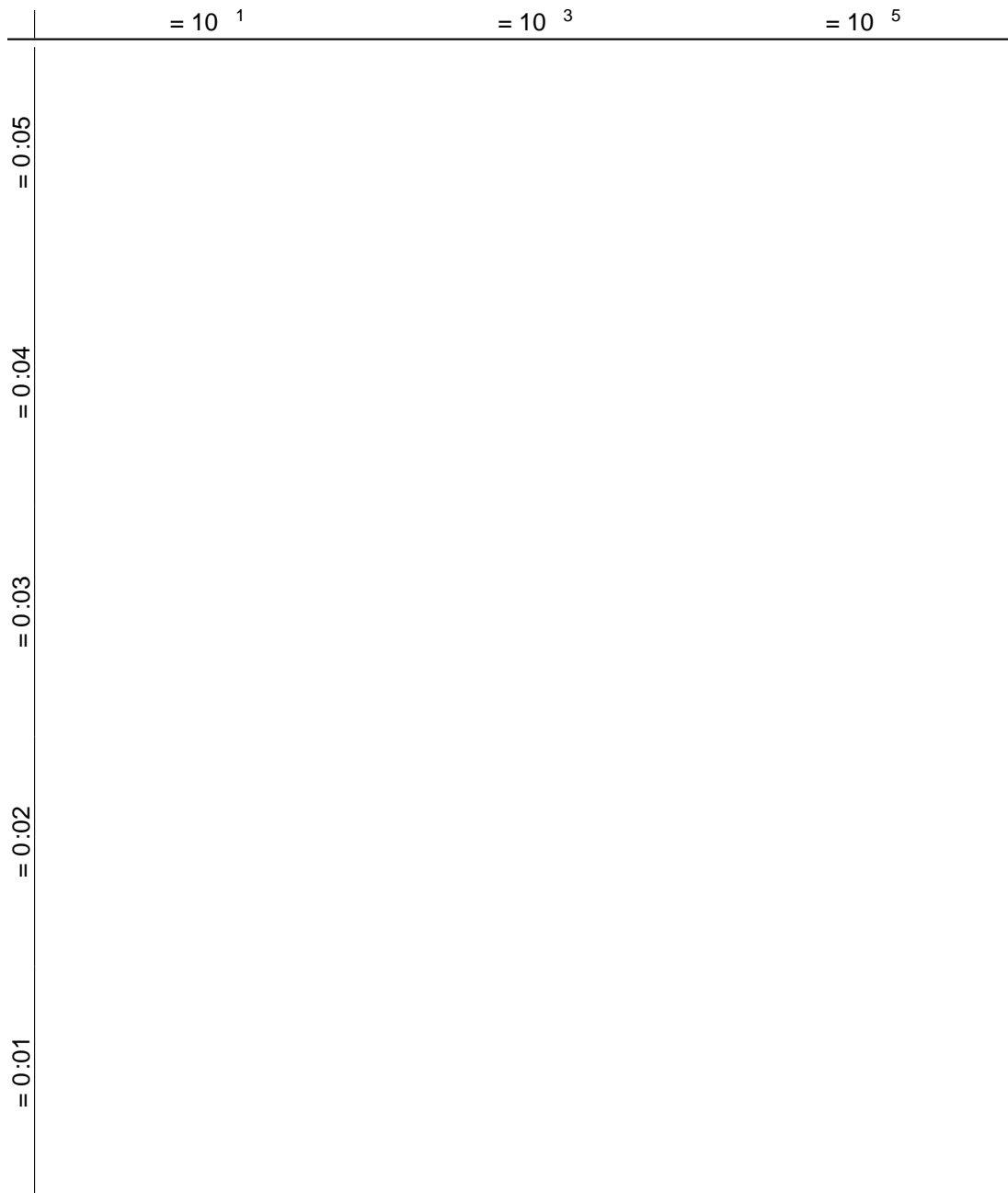


Table 7: Confidence set sizes for a neural network dynamics model trained on the half-cheetah environment, for varying ϵ and δ for $n = 5000$. The plots are as in Figure 3 (a).

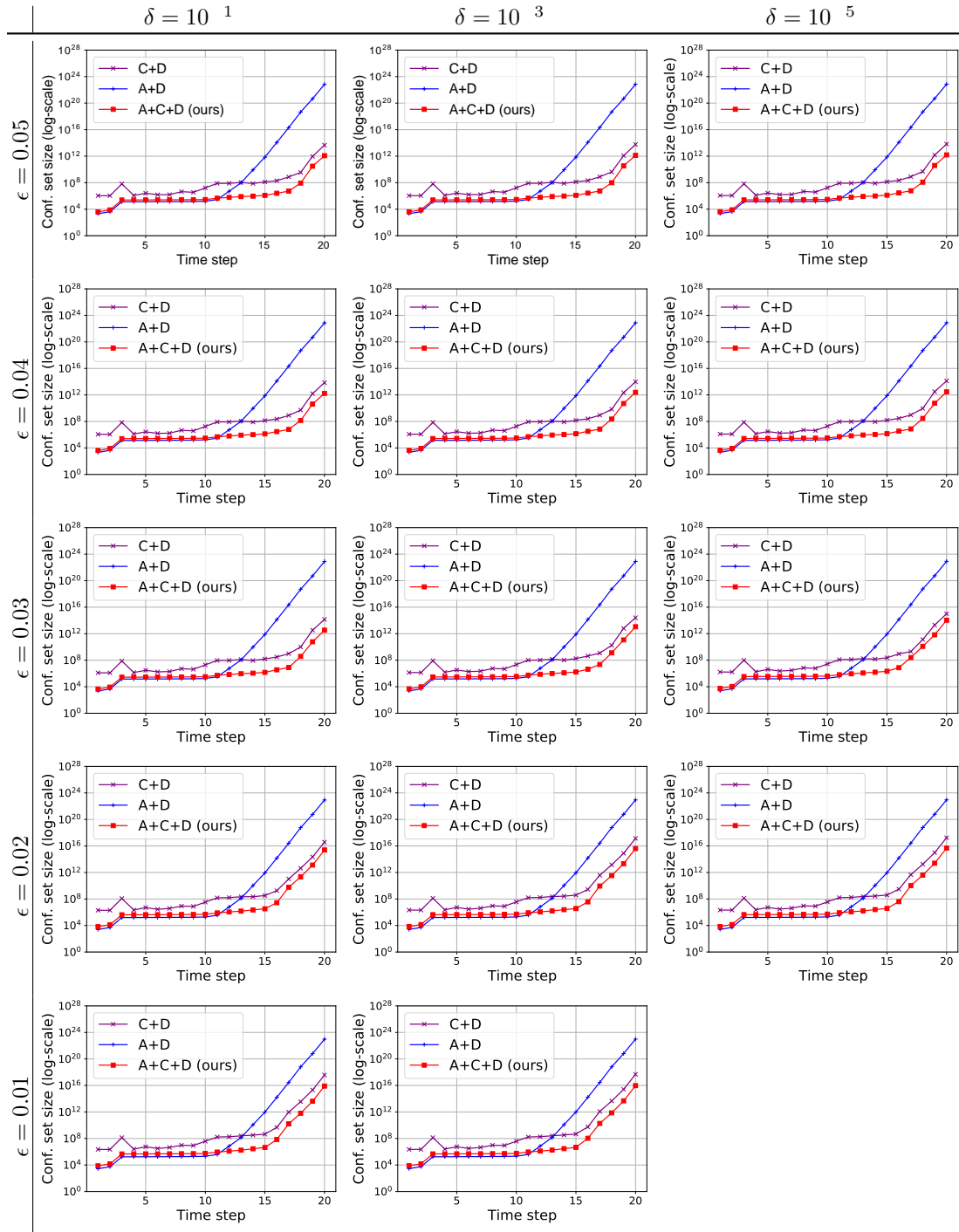


Table 8: Confidence set sizes for a neural network dynamics model trained on the half-cheetah environment, for varying ϵ, δ and for $n = 5000$. The plots are as in Figure 3 (b).

