

# Learning Prediction Intervals for Model Performance

Benjamin Elder,<sup>1</sup> Matthew Arnold,<sup>1</sup> Anupama Murthi<sup>1</sup> Jiri Navratil<sup>1</sup>

<sup>1</sup> IBM Research

benjamin.elder@ibm.com, marnold@us.ibm.com, anupama.murthi@ibm.com, jiri@us.ibm.com

## Abstract

1 Understanding model performance on unlabeled data is a fun- 41  
2 damental challenge of developing, deploying, and maintain- 42  
3 ing AI systems. Model performance is typically evaluated 43  
4 using test sets or periodic manual quality assessments, both 44  
5 of which require laborious manual data labeling. Automated 45  
6 *performance prediction* techniques aim to mitigate this bur- 46  
7 den, but potential inaccuracy and a lack of trust in their pre- 47  
8 dictions has prevented their widespread adoption. We address 48  
9 this core problem of performance prediction uncertainty with 49  
10 a method to compute *prediction intervals* for model perfor- 50  
11 mance. Our methodology uses transfer learning to train an *un-* 51  
12 *certainty model* to estimate the uncertainty of model perfor- 52  
13 mance predictions. We evaluate our approach across a wide 53  
14 range of drift conditions and show substantial improvement 54  
15 over competitive baselines. We believe this result makes pre- 55  
16 diction intervals, and performance prediction in general, sig- 56  
17 nificantly more practical for real-world use. 57

## 1 Introduction

18  
19 Knowing when a model’s predictions can be trusted is one of 58  
20 the key challenges in AI today. From an operational perspec- 59  
21 tive, understanding the quality of model predictions impacts 60  
22 nearly all stages in the model lifecycle, including pre-deploy 61  
23 testing, deployment, and production monitoring. From a so- 62  
24 cial perspective, prediction trust impacts society’s willing- 63  
25 ness to accept AI as it continues replacing human decision 64  
26 making in increasingly important roles. 65

27 Techniques such as *performance prediction* (Guerra, 66  
28 Prudêncio, and Ludermir 2008; Schat et al. 2020; Talagala, 67  
29 Li, and Kang 2019) strive to automatically predict the per- 68  
30 formance of a model with no human intervention. Unfortu- 69  
31 nately these techniques have not yet gained mainstream 70  
32 adoption, ironically enough, due to their potential unreli- 71  
33 ability and the resulting lack of trust in their predictions. Per- 72  
34 formance predictors are often surprisingly accurate when a 73  
35 base model is predicting on data similar to what it has al- 74  
36 ready seen in training and test. However, it is well known 75  
37 that model behavior can be extremely difficult to predict on 76  
38 previously unseen data (Nguyen, Yosinski, and Clune 2015; 77  
39 Su, Vargas, and Sakurai 2019). It is not reasonable to ex- 78  
40 pect a performance prediction algorithm to perfectly predict 79

a base model’s behavior in these scenarios. It is, however, 41  
reasonable to ask a performance predictor to quantify the 42  
uncertainty of its predictions so the application (or end user) 43  
can take appropriate precautions. 44

This paper describes a technique for computing *predic-* 45  
*tion intervals* on *meta-model* based performance predic- 46  
tions, to convey the degree to which the performance pre- 47  
diction should be trusted. Our technique uses *meta-meta-* 48  
*modeling* in a multi-task setting to train an *uncertainty model* 49  
to compute prediction intervals. Our approach makes no as- 50  
sumptions about the base model or performance predictor, 51  
and can easily be applied in other settings. The use of a sep- 52  
arate *meta-meta-model* to perform the uncertainty quantifi- 53  
cation allows the simultaneous prediction of both aleatoric 54  
(data-driven) and epistemic (model-driven) uncertainty. 55

A key challenge for an uncertainty model is predicting 56  
the unseen - ie, data that is substantially different than any- 57  
thing the model has seen in train or test. A simple cross val- 58  
idation or leave-one-out training is unlikely to produce suf- 59  
ficient feature-space drift to be informative in this regard. 60  
Our technique trains for these extreme cases by (1) simulat- 61  
ing various levels of drift ranging from mild to extreme, and 62  
(2) training on these drift scenarios using *external* data sets 63  
that are different from the base model training set but of the 64  
same modality. This approach enables the uncertainty model 65  
to learn how the performance predictor behaves in a variety 66  
of challenging scenarios, and uses this information to help 67  
predict risk when challenging scenarios arise in production. 68  
Due to its prevalence and commercial importance, we focus 69  
on tabular data, and therefore chose logistic regression and 70  
random forest base models. 71

We evaluate our uncertainty model on four different per- 72  
formance predictors and compare the uncertainty model 73  
against four different model-free baseline algorithms. In ev- 74  
ery scenario, our uncertainty model outperforms all base- 75  
lines, often by a large margin. Even without the use of an 76  
uncertainty model, our approach of using drift simulation 77  
for calibration yielded significant improvements over tradi- 78  
tional baselines. 79

## 2 Related work

80  
81 Prior work exists on performance prediction (Guerra, 82  
83 Prudêncio, and Ludermir 2008; Chen et al. 2019; Finn et al. 84  
2019; Redyuk et al. 2019; Schat et al. 2020; Talagala, Li, and 85

Kang 2019), however, to the best of our knowledge, none that assigns uncertainty bounds to their predictions. The field of domain generalization also includes some related work, for which a good recent review can be found in (Hospedales et al. 2020). For example, (Li et al. 2018) use a cross-domain meta-learning approach to model training similar our procedure, but applied to the base classification problem.

Of relevance are numerous methods for estimating the uncertainty of predictions from machine learning models in general. Many of these, ranging from classical statistical methods to state-of-the-art deep learning (DL) models ((Gal and Ghahramani 2016; Kendall and Gal 2017; Koenker and Bassett 1978; Nix and Weigend 1994), also see (Khosravi et al. 2011) for a review), could be applied to meta-model based performance prediction.

There are well-established parametric methods for prediction intervals, see for example (Geisser 2017). Methods implicitly learning the error distribution are also available, for example by incorporating a feature-dependent variance into the loss function for iterative training procedures (Nix and Weigend 1994). Furthermore, there has been significant progress in constructing neural architectures which simultaneously output a classification and an uncertainty prediction (Brosse et al. 2020; Kabir et al. 2019; Khosravi, Nahavandi, and Creighton 2010; Malinin and Gales 2018).

Non-parametric methods such as the jackknife and bootstrap (Efron 1979; Efron and Gong 1983), and more recent variations (Lei et al. 2018; Vovk et al. 2018; Papadopoulos 2008; Vovk 2012; Vovk, Gammerman, and Shafer 2005) can estimate the uncertainty of a statistical prediction without assuming a particular error model, but rely on the assumption that the distribution of the unlabeled data is the same as the train data. Ensemble methods have been proposed for both traditional ML (Dietterich 2000; Kwok and Carter 1990) and DL models (Hansen and Salamon 1990; Lakshminarayanan, Pritzel, and Blundell 2017; Osband, Aslanides, and Cassirer 2018). The variance of the ensemble predictions can be used as a feature-dependent measure of uncertainty. We implement this strategy as one of our baselines (see Sec. 4.3). Bayesian approaches have been extended to non-parametric applications, including neural networks (Bishop 1997; Blundell et al. 2015; Neal et al. 2011) and a popular approximate version of Bayesian neural networks - the Monte-Carlo dropout approximation (Gal and Ghahramani 2015; Gal, Hron, and Kendall 2017) - also serves as a baseline in our work.

### 3 Method

Our approach estimates prediction intervals for the performance of a black-box classification model on a pool of unlabeled data. First, we use meta-modeling to predict the accuracy of the base classification model (*performance prediction*). Second, a pre-trained *uncertainty model* is used to estimate a prediction interval, which describes the probable range for the true value of the accuracy.

#### 3.1 Performance predictors

As its name suggests, *performance prediction* is the problem of estimating the value of a performance metric that a

machine-learning model will achieve for a given pool of unlabeled data (referred to here as the production set). This work focuses on classification accuracy, but the same methods could be applied to other metrics such as the F1 score or the error of a regression model. We treat the base model as a black box from which only the vector of predicted class probabilities for each sample is available. We developed two types of performance predictors (four variants in all) which we will use as the basis for our prediction intervals.

The performance predictors that we use in this work each output a confidence score for each unlabeled data point in the production set. This score, between zero and one, is an estimate of the likelihood that the base model predicted the correct class label. For the purposes of computing an aggregate accuracy score for the production set, we take the average of these confidences. This confidence averaging produced better estimates of the accuracy than making binary correct/incorrect predictions for each sample.

The `confidence` predictor is a simple, binning-based procedure, which recalibrates the base model confidence score for the most likely class (Zadrozny and Elkan 2001). The values of this confidence on the test set are gathered into a histogram (binned in increments of 0.1), and the base model accuracy is computed for each bin. A performance prediction for a data point is given by the average accuracy of the bin that spans that point's base model confidence.

Our `meta-model` performance predictor uses its own model to predict data points that are likely to be mislabeled by the base model. Training data for the `meta-model` predictor is created by relabeling the test set with binary labels indicating whether the base model correctly classified each sample. The meta-model, which is an ensemble of a Gradient Boosting Machine (GBM) and a logistic regression model, classifies each sample as correct or incorrect, and the probability assigned to the "correct" class is returned as the performance predictor confidence score. Further details of the `meta-model` predictor implementation are provided in the supplementary material.

#### 3.2 Uncertainty model

We propose a new approach for computing prediction intervals by using an *uncertainty model* (a meta-meta-model) that learns to predict the behavior of a performance predictor (a meta-model). This uncertainty model (UM) is a regression model that quantifies the uncertainty of the performance predictor's estimate of the base model accuracy on the production set. We pre-train the UM using a library of training datasets. The UM can then observe the behavior of a performance predictor on a new target dataset and generate a prediction interval.

The UM must be trained using examples of performance prediction errors. Each training sample for the UM consists of a full drift scenario, comprising: (1) labeled train and test datasets, (2) a pool of unlabeled data (the production set), (3) a base classification model trained on the train set, and (4) a performance predictor trained for this dataset and base model. We generate a large number of such examples using the simulation procedures described in Sec. 4.1. The target values are the (absolute) differences between the true and

199 predicted accuracy on the entire production set. This method  
200 could be extended to predict the signed value of the errors,  
201 allowing for asymmetric prediction intervals.

202 The UM architecture is an ensemble of GBM models. Ex-  
203 perimentally, we found that an ensemble of ten models re-  
204 duced prediction variance and led to improved accuracy. The  
205 model was trained using a quantile loss function, which nat-  
206 urally enables the calculation of prediction intervals targeted  
207 to capture the true error with a specified probability. Further  
208 implementation details are described in the supplementary  
209 material.

210 **Features** The UM was trained using a set of derived  
211 features that are generic enough to be compatible across  
212 datasets with varying feature spaces and numbers of classes.  
213 The derived features are extracted from a number of mod-  
214 els, including (1) the base classification model, (2) the per-  
215 formance predictor, (3) a group of proxy models, and (4) a  
216 group of drift models. The proxy models (one logistic re-  
217 gression, one random forest, and one GBM) were trained on  
218 the same features and classification task as the base model,  
219 and provide a complementary perspective on the classifica-  
220 tion difficulty of each data point. The drift models are ran-  
221 dom forest models trained to predict whether a given sample  
222 came from the test set or the production set. They provide di-  
223 rect insight into the degree of feature space drift for a given  
224 scenario. Further implementation details for the proxy and  
225 drift models are provided in the supplementary material.

226 A full list of the features used for the UM is shown in  
227 Table 1. The procedure for constructing the features of type  
228 *Distance* starts by choosing a function  $f$  that maps any  
229 feature vector to a scalar value, for example the highest  
230 value from the base model confidence vector. The value  
231 of this function is used to construct two histograms, one  
232 for the samples in the test set, and one for the production  
233 set. Finally, the distance between these two histograms is  
234 computed using some distance function  $D$ . We used three  
235 functions for  $D$ : the Kolmogorov–Smirnov metric  $D_1 =$   
236  $\max_i |P_i - Q_i|$ , the inverse overlap  $D_2 = \sum_i \max(p_i -$   
237  $q_i, 0)$ , and the squared inverse overlap  $D_3 = \sum_i (\max(p_i -$   
238  $q_i, 0))^2$ . Here  $p$  and  $q$  are the normalized histograms from  
239 the test and production sets,  $i$  indexes the bins in the his-  
240 tograms (which must be identically spaced), and  $P, Q$  are  
241 the CDFs corresponding to  $p, q$ .

242 The remaining features can be grouped into three  
243 types: *Prediction*, *Noise*, and *Internal*. The  
244 *Prediction* features are directly derived from the pre-  
245 dictions of one of the source models (without using the  
246 *Distance* procedure). These include the entropy of the  
247 base model predictions, the change in accuracy predicted by  
248 the performance predictor, and the accuracy of the drift clas-  
249 sifiers. The *Noise* features are the size of bootstrap confi-  
250 dence intervals for the average of the top base model confi-  
251 dence and the performance predictor confidence, both of  
252 which approach zero as the number of points in the pro-  
253 duction set approaches infinity. The *Internal* features are  
254 white-box quantities extracted from the performance predic-  
255 tor or proxy models, such as the *intrinsic* prediction inter-

256 vals described in Sec. 4.3, or the difference between the cal-  
257 ibrated and uncalibrated performance predictions. A com-  
258 plete description of the features listed here is given in the  
259 supplementary material.

260 An ablation study is presented in Table 2, showing the  
261 performance of the UM, coupled with the *meta-model*  
262 performance predictor, obtained using subsets of the full  
263 set of features described above. This study shows that the  
264 most effective information for understanding the perfor-  
265 mance prediction uncertainty comes from the drift models.  
266 It also shows that including the more numerous and compu-  
267 tationally expensive *Distance* and *Internal* features  
268 deliver a significant performance boost over the simpler  
269 *Prediction* and *Noise* features.

## 4 Experimental methodology 270

271 Fig. 1 shows example results for *linear-skew* drift sce-  
272 narios (Sec. 4.1) simulated from the *bng-zoo* dataset. The  
273 middle and right panel show accuracy predictions from  
274 the *meta-model* performance predictor, and the uncertainty  
275 model prediction intervals calibrated using two different lev-  
276 els of  $\alpha$ , (0.5 and 0.9), as described in Sec. 4.5.

### 4.1 Drift simulation 277

278 Training and evaluating the UM requires examples of data  
279 drift. The breadth of the drift examples used for training  
280 largely determines the quality and coverage of the result-  
281 ing model. Tabular labeled datasets containing sufficient nat-  
282 urally occurring drift are difficult to obtain, therefore we  
283 chose to generate such examples through resampling-based  
284 simulation. We focus on *covariate shift* because it has been  
285 shown to encompass a wide range of real-world drift scenar-  
286 ios (Card and Smith 2018).

287 We use two different algorithms for generating drift: (1)  
288 *linear-skew*, is designed for breadth of coverage, for  
289 providing training data, and ensuring the generation of ex-  
290 treme drift (2) *nearest-neighbors* is designed to sim-  
291 ulate drift more likely to occur in the real world for an ad-  
292 ditional evaluation scenario. Examples of drift generated by  
293 both algorithms is included in the supplementary material.

294 **Linear-skew** The *linear-skew* method requires  
295 choosing a feature dimension ( $F$ ) along which the bias will  
296 be induced, and a threshold  $t$  to split the dataset into two  
297 buckets. The sampling parameter  $R$  controls the ratio of  
298 sampling from the two buckets for the train/test sets and the  
299 production set, thus also controlling the amount of drift in  
300 the scenario. When  $R = 50$ , the train, test, and production  
301 sets all have the same distribution, and there is no drift.  
302 When  $R = 0$  or  $R = 100$ , there is no overlap between the  
303 train/test distribution and the production distribution. The  
304 *linear-skew* procedure is described in Alg. 1.

305 **Nearest-neighbors** The *nearest-neighbors* algo-  
306 rithm strives to simulate a particular demographic either ap-  
307 pearing, or disappearing from production traffic. It does so  
308 by sampling a data point and then uses nearest neighbors to  
309 identify other data points that are similar (nearest neighbors)

Source	Name	Description	Type	Count
Base	top confidence	confidence score of predicted class	D	3
	(top-2nd) confidence	top - second highest class probability	D	3
	confidence entropy	entropy of confidence vector	D	3
	class frequency	relative frequency of predicted classes	D	3
	entropy ratio	avg. prod. set entropy/avg. test set entropy	P	1
	bootstrap	size of bootstrap conf. intervals for avg. accuracy	N	1
Perf. Pred.	predicted change	predicted prod. acc. - base test acc.	P	1
	avg. pred. stdev.	stdev. of confidence scores (prod)	P	1
	avg. pred. entropy	entropy of confidence histogram (prod)	P	1
	predicted uncertainty	intrinsic uncertainty interval	I	1
	bootstrap	size of bootstrap conf. intervals for pred. acc.	N	1
	ensemble whitebox	whitebox stats from meta-model ensemble	I	2
	calibration whitebox	whitebox stats from calibration	I	2
Proxy	gbm whitebox	internal stats from gbm meta-model	I	16
	top confidence	confidence score of predicted class	D	9
	(top-2nd) confidence	top - second highest class probability	D	9
	best feature	projection onto most important feature	D	3
Drift	num import. feat.	num. features to make 90% feat. importance	I	1
	accuracy	test vs. prod classification accuracy	P	3
Other	(top-2nd) confidence	top - second highest class probability	D	9
	PCA projection	projection onto highest PCA component	D	3

Table 1: Source models, names, descriptions, types (D=Distance, P=Prediction, N=Noise, I=Internal), and counts for all UM features.

	Average Cost ( $\alpha = 0.9$ )				
	Base	Perf. Pred.	Proxy	Drift	All
Distance	1.46	–	1.21	1.18	1.07
Internal	1.61	1.15	–	–	1.05
Prediction	–	1.29	1.87	–	1.33
Noise	1.46	1.59	–	–	1.32
All	1.30	1.25	1.21	1.16	1.0

Table 2: Average uncertainty model cost, Eq. (1), for the meta-model predictor, with  $\alpha = 0.9$ , normalized by the value when all features are used.

#### Algorithm 1 Algorithm to create linear-skew drift scenarios

**Input:** Dataset  $X \subset \mathcal{X}$ ;  $p_{tr}, p_{te}, p_{pr} \in [0, 1] : p_{tr} + p_{te} + p_{pr} = 1$ ; Feature dimension  $F : \mathcal{X}^{(F)} \subset \mathcal{X}$ ; threshold function  $t : \mathcal{X}^{(F)} \rightarrow \{X_A, X_B\}$ ; Sampling ratio  $R \in [0, 100]$ ; minibatch size  $b$

**Output:**  $X_{tr}, X_{te}, X_{pr} \subset X$   
 $X_A, X_B, X_{tt}, X_{pr} \leftarrow \{\}$

**for**  $x$  **in**  $X$  **do**  
  Add  $x$  to  $t(X^{(F)}) \in \{X_A, X_B\}$  ▷ Add data point to bucket, defined by threshold  $t$   
**end for**  
**while**  $|X_A| > b$  **and**  $|X_B| > b$  **do** ▷ Randomly sample  $X_{tt}, X_{pr}$  from buckets until out of data points  
  Add  $(p_{tr} + p_{te}) \times \frac{R}{100} \times b$  points from  $X_A$  and  $(p_{tr} + p_{te}) \times (1 - \frac{R}{100}) \times b$  points from  $X_B$  into  $X_{tt}$   
  Add  $p_{pr} \times (1 - \frac{R}{100}) \times b$  points from  $X_A$  and  $p_{pr} \times \frac{R}{100} \times b$  points from  $X_B$  into  $X_{pr}$   
**end while**  
  Randomly split  $X_{tr}, X_{te} \leftarrow X_{tt}$  with proportions  $\frac{p_{tr}}{p_{tr} + p_{te}}$  and  $\frac{p_{te}}{p_{tr} + p_{te}}$

maining training datasets. All results are averaged over these 330  
fifteen different UMs. 331

Since we focused on tabular data, we used random 332  
forest and logistic regression base models for all ex- 333  
periments. In the linear-skew simulations we chose 334  
two features per dataset, and performed Alg. 1 for 335  
each feature with fifteen values of the sampling ratio 336  
 $R = 0, 1, 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 95, 99, 100$ . 337  
This was repeated using five random seeds, giving 338

310 or dissimilar (furthest neighbors) and remove them from  
311 the dataset. We used this algorithm to create fairly severe  
312 drift, removing 50-70% of the original data points, which  
313 we believe covers the range of realistic possible drift. The  
314 nearest-neighbors algorithm is described in Alg. 2.

## 4.2 Datasets and settings

316 For our experiments we use a set of fifteen publicly avail-  
317 able tabular datasets, sourced from Kaggle, OpenML, and  
318 Lending Club: Artificial Character, Bach Choral, Bank Mar-  
319 keting, BNG Zoo, BNG Ionosphere, Churn Modeling, Cred-  
320 itcard Default, Forest Cover Type, Higgs Boson, Lending  
321 Club (2016 Q1, 2017 Q1), Network Attack, Phishing, Pul-  
322 sar, SDSS, and Waveform. Details of the individual dataset’s  
323 characteristics and our pre-processing procedures are pro-  
324 vided in the supplementary material.

325 To simulate the effect of training the UM on an offline  
326 library of training datasets and then deploying it to make  
327 predictions on a new, unseen target dataset, we conducted  
328 our experiments in a leave-one-out manner. Each dataset was  
329 chosen in turn as the target, and its UM was trained on the re-

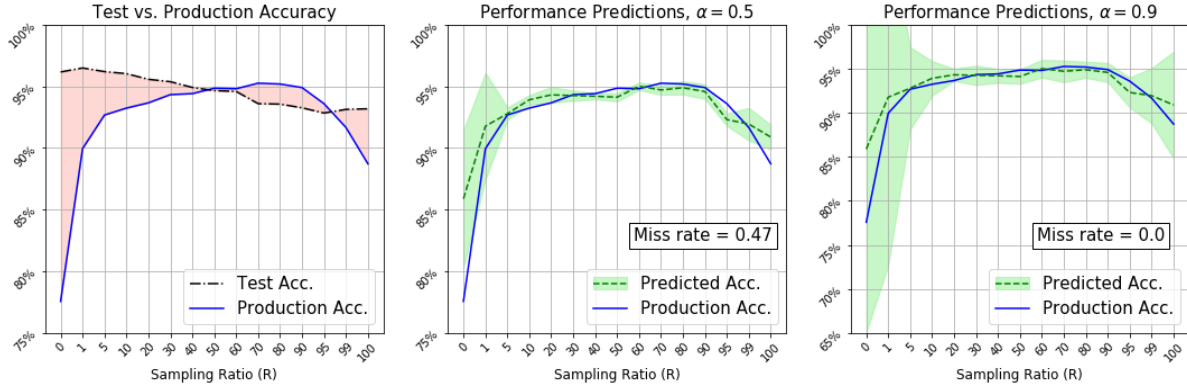


Figure 1: Example of performance prediction and the UM for linear-skew drift scenarios. The left plot shows the model accuracy drift (shaded area) induced by Alg. 1. The middle (right) plot shows the accuracy predicted by the meta-model predictor, with UM prediction intervals calibrated to  $\alpha = 0.5$  ( $\alpha = 0.9$ ).

---

**Algorithm 2** Algorithm to create nearest-neighbors drift scenarios

---

**Input:** Dataset  $X \subset \mathcal{X}$ ;  $p_{tr}, p_{te}, p_{pr} \in [0, 1] : p_{tr} + p_{te} + p_{pr} = 1$ ;  $P_{set} \in [0, 1]$ ;  $P_{near} \in [0, 1]$ ;  $P_{down} \in [0, 1]$   
**Output:**  $X_{tr}, X_{te}, X_{pr} \subset X$   
 Randomly split data into  $X_{pr}$  and  $X_{tt}$  with proportions  $p_{pr}$  and  $1 - p_{pr}$   
 With probability  $P_{set}$ , set downsample set  $X_{down} = X_{tt}$  and  $X_{rand} = X_{pr}$ , else  $X_{down} = X_{pr}$  and  $X_{rand} = X_{tt}$   $\triangleright$   
 Choose distribution to bias non-randomly  
 Choose point  $p \in X_{down}$  at random  
 Order points  $x \neq p \in X_{down}$  by distance from  $p$   
 Choose  $D$ =nearest ( $N$ ) with probability  $P_{near}$  else  $D$ =furthest ( $F$ )  $\triangleright$  Choose nearest or furthest bias  
 Remove the fraction  $P_{down}$  points which are  $D \in \{N, F\}$  from  $p$   
 Remove fraction  $P_{down}$  from  $X_{rand}$  randomly  $\triangleright$  Randomly downsample non-biased distribution  
 Randomly split  $X_{tr}, X_{te} \leftarrow X_{tt}$  with proportions  $\frac{p_{tr}}{p_{tr} + p_{te}}$  and  $\frac{p_{te}}{p_{tr} + p_{te}}$

---

compute an uncertainty score for each point in the  $k$ -th bin as  $u_k = \sqrt{a_k(1 - a_k)/n_k}$ , which is the standard error of a Bernoulli distribution with parameter  $a_k$ .

For the meta-model predictor, we created two variants that replace the GBM and logistic regression meta-models with different classifiers. The `crossval` predictor uses an ensemble of ten random forest models, each trained with a different cross-validation fold of the test set, and the standard deviation of their predictions is used as the uncertainty estimate. The `dropout` predictor uses an XGBoost (Chen and Guestrin 2016) model with the DART (Vinayak and Gilad-Bachrach 2015) booster, which applies dropout (Srivastava et al. 2014) regularization to GBM models. In the spirit of the Monte-Carlo dropout approach for Bayesian neural networks (Gal and Ghahramani 2015; Gal, Hron, and Kendall 2017), ten predictions are made for each sample with dropout turned on to introduce randomness in the confidence scores, giving an estimated average and standard deviation for the model accuracy.

Besides the intrinsic baselines, we also compare with three other baseline methods, which produce prediction intervals using: (1) the standard error of the mean of the performance predictor confidences, (2) the size of a bootstrap uncertainty interval for the mean of these confidences, and (3) a constant sized prediction interval.

#### 4.4 Evaluation Metric

Evaluating the quality of a set of prediction intervals involves a trade-off between two opposing kinds of errors: prediction intervals that are too small and do not capture the magnitude of the true error (Type I cost), and prediction intervals that are unnecessarily large (Type II cost). In an ad hoc comparison between two methods generating prediction intervals, it is common that one method does not dominate the other in the sense of having both smaller Type I and Type II error. A comprehensive comparison of such methods requires making a tradeoff between the two.

One common approach to quantifying this trade-off is to scale each set of prediction intervals by a constant factor to

339 a total of 300 drift scenarios per dataset.<sup>1</sup> For the  
 340 nearest-neighbors simulations, 300 scenarios were  
 341 generated for each dataset/base model combination with pa-  
 342 rameters  $P_{set} = 0.5$ ,  $P_{near} = 0.5$ , and  $P_{down} \in [0.5, 0.7]$ ,  
 343 for a total of 9000 scenarios.

#### 4.3 Model-Free Baselines

345 A set of model-free baseline techniques for computing pre-  
 346 diction intervals are compared against the UM. The first  
 347 set of techniques are three *intrinsic* methods which lever-  
 348 age white-box information from the performance predictors.  
 349 These intrinsic methods produce uncertainty estimates for  
 350 each point in the production set, and the average of these  
 351 estimates is used as the (uncalibrated) prediction interval.

352 For the `confidence` predictor, if the accuracy in bin  $k$   
 353 is  $a_k$  and the number of samples falling into bin  $k$  is  $n_k$ , we

<sup>1</sup>Except for the Network Attack dataset, which only has one feature amenable to this procedure, see supplementary material.

achieve a common miss rate (eg 5%), and then compare their average size or average excess. We chose instead to evaluate results based on a cost function which penalizes both types of error:

$$C_\alpha(\vec{\delta}, \vec{u}) = \sum_i \left[ \alpha \max(\delta_i - u_i, 0) + (1 - \alpha) \max(u_i - \delta_i, 0) \right]. \quad (1)$$

In Eq. (1),  $\delta_i = |a_i - p_i|$  is the difference between the true base model accuracy and the performance prediction for production set  $i$ , and  $u_i$  is the (single-directional) size of the scaled prediction interval. The cost function parameter  $\alpha \in [0, 1]$  can be adjusted to control the balance between the two types of error.

#### 4.5 Calibration

Prediction intervals must be calibrated in order to achieve reliable performance for a scale-sensitive metric such as Eq. (1). We compare two calibration methods, one with external drift scenarios as a holdout set, and one without. Without the holdout set for calibration (Target Calibration), we used the approximation that the performance prediction error is normally distributed, and multiplied the uncalibrated prediction intervals by the Z-score of the desired confidence interval (determined by the cost function parameter  $\alpha$ ).<sup>2</sup> This calibration method is applied to all of the model-free baselines described in Sec. 4.3.

The second calibration method (TL Calibration) used the simulated drift scenarios from external holdout datasets as a calibration set. A constant scale factor was computed which minimized the cost function Eq. (1) for each set of prediction intervals on this hold-out set. This calibration method was applied to the model-free baselines using 100% of the holdout drift scenarios, and to the UM prediction intervals using an 80%/20% train/test split of the drift scenarios.

### 5 Experimental results

In this section we demonstrate the performance of our UM, trained using the `linear-skew` drift scenarios, and evaluated using both the `linear-skew` and the `nearest-neighbors` style scenarios. We use the four performance predictors described in Sec. 3.1, and the four baselines described in Sec. 4.3. The quality of the prediction intervals is measured by the cost function  $C_\alpha$  from Eq. (1), with  $\alpha$  between 0.5 and 0.95. This range covers most reasonable user preferences for penalizing under- vs. over-shooting of the appropriate prediction interval size.

**Linear-skew results** Fig. 2 compares the cost  $C_\alpha$  of the UM and the standard error, bootstrap, and intrinsic baselines described in Sec. 4.3, evaluated using the `linear-skew`

<sup>2</sup>The cost is minimized by balancing the two terms in Eq. (1), thus the  $\alpha$ -confidence interval chosen for calibration. For example, if  $\alpha = 0.95$  is chosen, then the scale factor is  $Z = 1.96$ .

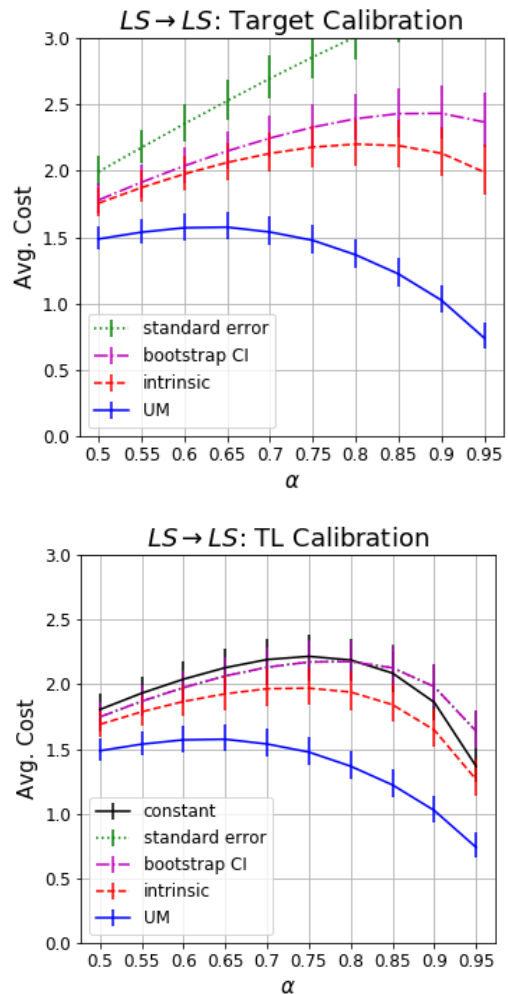


Figure 2: Average cost (Eq. (1)) for prediction intervals. The baselines from Sec. 4.3, calibrated without (top) and with (bottom) external datasets, are compared to the UM. The evaluation uses the `linear-skew` drift scenarios.

drift scenarios. The results are averaged across the four performance predictors<sup>3</sup>, and calibrated using the target dataset method (top) and external holdout drift scenarios (bottom), as described in Sec. 4.5. The error bars in Fig. 2 indicate the 95% bootstrap confidence interval. The bottom plot also includes the calibrated constant baseline.

It is clear from the upper panel of Fig. 2 that the UM method trained with the `linear-skew` simulated drift scenarios substantially outperforms the baselines. This is especially true for moderate to high values of  $\alpha$ , which correspond to penalizing prediction intervals that are too small more than intervals which are too large. Comparing this result to the bottom panel, we see that merely calibrating with the simulated drift scenarios can dramatically improve the performance of the baseline methods. However, the UM still

<sup>3</sup>There are only three sets of results for the “intrinsic” curve, as the `meta-model` predictor has no intrinsic prediction interval.

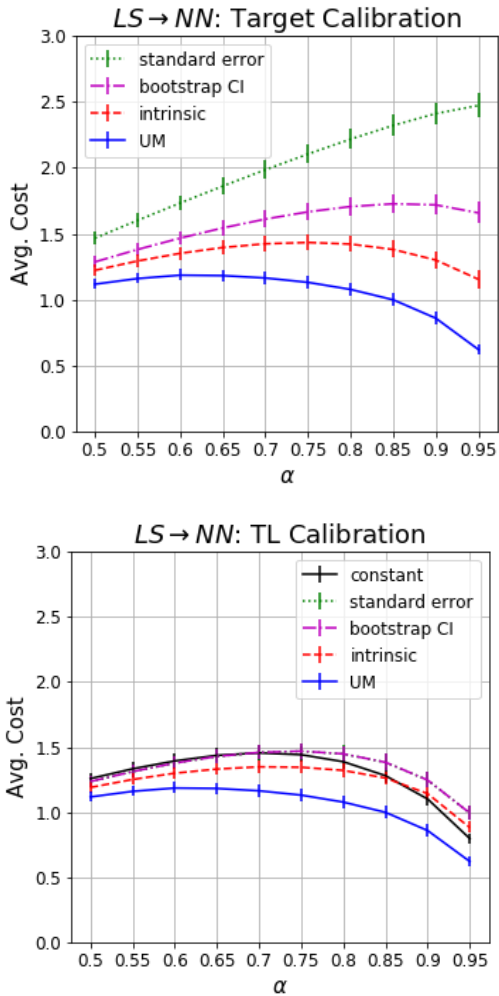


Figure 3: Same experiments as Fig. 2, except using nearest-neighbors drift scenarios for evaluation.

445 provides a major improvement at all values of  $\alpha$ .

446 **Nearest-neighbors results** Fig. 3 shows the same experi- 447  
 448 ments as Fig. 2, but using nearest-neighbors scenar- 449  
 450 ios for evaluation (and still using linear-skew scenar- 451  
 452 ios for training). The UM still outperforms the baselines for the 453  
 454 full range of  $\alpha$ . The overall costs tend to be smaller, since the 455  
 456 model accuracy drift in the nearest-neighbors scenar- 457  
 458 ios was smaller on average than in the linear-skew 459  
 460 scenarios. This result confirms that the deliberately en- 461  
 462 gineered linear-skew scenarios are able to provide 463  
 464 effective training for more realistic, organically created 465  
 466 nearest-neighbors drift scenarios.

467 Table 3 provides all of the results in numeric form, bro- 468  
 469 ken down by performance predictor, and averaged across the 470  
 471 same range of  $\alpha$  values. This confirms that the individual 472  
 473 predictor results align with the averaged results previously 474  
 475 shown. 476  
 477  
 478

LS (train) $\rightarrow$ LS (eval)				
Predictor				
Method	Confidence	Crossval	Dropout	Meta-model
SE	2.88	2.57	3.09	2.41
BS	2.38	2.06	2.50	1.89
I	2.23	1.72	2.20	–
SE (TL)	2.13	1.91	2.13	1.77
BS (TL)	1.13	1.91	2.13	1.77
C (TL)	2.11	1.88	2.10	1.83
I (TL)	1.83	1.67	1.86	–
UM	<b>1.46</b>	<b>1.24</b>	<b>1.38</b>	<b>1.33</b>
LS (train) $\rightarrow$ NN (eval)				
Predictor				
Method	Confidence	Crossval	Dropout	Meta-model
SE	1.97	1.87	2.48	1.75
BS	1.59	1.41	2.00	1.29
I	1.39	1.14	1.48	–
SE (TL)	1.36	1.25	1.58	1.15
BS (TL)	1.36	1.25	1.58	1.15
C (TL)	1.30	1.20	1.49	1.16
I (TL)	1.26	1.17	1.29	–
UM	<b>1.08</b>	<b>1.00</b>	<b>1.14</b>	<b>0.97</b>

Table 3: Costs for both experiments, including the standard error (SE), bootstrap (BS), and intrinsic (I) baselines, the same methods with our transfer-learning (TL) calibration, as well as the calibrated constant (C) and UM, averaged over the values of  $\alpha$  shown in the figures (0.5 to 0.95).

462 **Generalization and Limitations** Choosing the source 463  
 464 datasets for this transfer-learning based approach is an im- 465  
 466 portant consideration. For a domain specific application, it 467  
 468 is obviously preferable to choose source datasets from the 469  
 470 same or a closely related domain. In addition to the domain, 471  
 472 we expect that it is valuable to approximately match other 473  
 474 dataset characteristics such as number of classes, number of 475  
 476 features, feature sparsity, etc. 477  
 478

479 For applications to other data modalities, for example im- 480  
 481 age or text data, many standard benchmark datasets such as 482  
 483 ImageNet could be used for pre-training. The models used in 484  
 485 the UM and the feature extraction would need to be replaced 486  
 487 or supplemented with modality appropriate architectures to 488  
 489 extract meaningful features from the data. Finally, the base 490  
 491 models used in the drift scenarios for training the UM should 492  
 493 include models of the same class as those to which it will be 494  
 495 applied at prediction time. 496  
 497  
 498

## 6 Conclusion 479

480 Performance prediction is an invaluable part of the deploy- 481  
 482 ing, monitoring, and improving an AI model. This paper 483  
 484 addresses this problem by describing a novel technique for 485  
 486 quantifying model uncertainty. It leverages multi-task learn- 487  
 488 ing and meta-meta-modeling to generate prediction inter- 489  
 490 vals on any model-based performance prediction system. 491  
 492 Our method substantially outperforms a group of competi- 493  
 494 tive baselines on dataset shift produced by two different sim- 495  
 496 ulation mechanisms. We believe this work helps make per- 497  
 498 formance prediction more practical for real-world use, and 499  
 499 may encourage further innovation in this important area. 500

## 7 Ethics Statement

Training AI models to “know what they don’t know” is one of the key challenges in AI today (Kindig 2020; Davey 2018; Knight 2018). AI models can be notoriously overconfident in scenarios that their training data did not prepare them for (Nguyen, Yosinski, and Clune 2015), eroding trust in AI and society’s willingness to accept AI as it continues to replace human decision making in increasingly important roles. Our work directly addresses this problem, proposing a novel technique for quantifying model uncertainty that beats all baselines we compare against. In addition, we hope that this work will encourage the community to continue to invest and innovate in this important area. Due to the nature of this problem we do not foresee any negative consequences stemming from our work.

## References

Bishop, C. M. 1997. Bayesian Neural Networks. *Journal of the Brazilian Computer Society* 4. ISSN 0104-6500.

Blundell, C.; Cornebise, J.; Kavukcuoglu, K.; and Wierstra, D. 2015. Weight Uncertainty in Neural Networks. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning*, volume 37 of *ICML’15*, 1613–1622. JMLR.org.

Brosse, N.; Riquelme, C.; Martin, A.; Gelly, S.; and Moulines, É. 2020. On Last-Layer Algorithms for Classification: Decoupling Representation from Uncertainty Estimation. *arXiv preprint arXiv:2001.08049*.

Card, D.; and Smith, N. A. 2018. The Importance of Calibration for Estimating Proportions from Annotations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, 1636–1646. Association for Computational Linguistics.

Chen, T.; and Guestrin, C. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’16, 785–794. Association for Computing Machinery. ISBN 9781450342322.

Chen, T.; Navratil, J.; Iyengar, V.; and Shanmugam, K. 2019. Confidence Scoring Using Whitebox Meta-models with Linear Classifier Probes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, 1467–1475.

Davey, T. 2018. How AI Handles Uncertainty. <https://futureoflife.org/2018/03/15/how-ai-handles-uncertainty-brian-ziebart>.

Dietterich, T. G. 2000. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, 1–15. Springer. ISBN 978-3-540-45014-6.

Efron, B. 1979. Bootstrap Methods: Another Look at the Jackknife. *The Annals of Statistics* 7(1): 1–26.

Efron, B.; and Gong, G. 1983. A Leisurely Look at the Bootstrap, the Jackknife, and Cross-Validation. *The American Statistician* 37(1): 36–48.

Finn, C.; Rajeswaran, A.; Kakade, S.; and Levine, S. 2019. Online Meta-Learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, 1920–1930. PMLR.

Gal, Y.; and Ghahramani, Z. 2015. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *Proceedings of The 33rd International Conference on Machine Learning*.

Gal, Y.; and Ghahramani, Z. 2016. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. In *Advances in Neural Information Processing Systems*, volume 29, 1019–1027. Curran Associates, Inc.

Gal, Y.; Hron, J.; and Kendall, A. 2017. Concrete Dropout. In *Advances in Neural Information Processing Systems*, volume 30, 3581–3590. Curran Associates, Inc.

Geisser, S. 2017. *Predictive Inference*. CRC Press. ISBN 9781351422291.

Guerra, S. B.; Prudêncio, R. B.; and Ludermir, T. B. 2008. Predicting the performance of learning algorithms using support vector machines as meta-regressors. In *International Conference on Artificial Neural Networks*, volume 5163, 523–532. Springer.

Hansen, L.; and Salamon, P. 1990. Neural Network Ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12: 993–1001. ISSN 0162-8828.

Hospedales, T.; Antoniou, A.; Micaelli, P.; and Storkey, A. 2020. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*.

Kabir, H.; Khosravi, A.; Kavousi-Fard, A.; Nahavandi, S.; and Srinivasan, D. 2019. Optimal Uncertainty-guided Neural Network Training. *arXiv preprint arXiv:1912.12761*.

Kendall, A.; and Gal, Y. 2017. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? In *Advances in Neural Information Processing Systems*, volume 30, 5574–5584. Curran Associates, Inc.

Khosravi, A.; Nahavandi, S.; and Creighton, D. 2010. Construction of Optimal Prediction Intervals for Load Forecasting Problems. *IEEE Transactions on Power Systems* 25(3): 1496–1503. ISSN 1558-0679.

Khosravi, A.; Nahavandi, S.; Creighton, D.; and Atiya, A. F. 2011. Comprehensive Review of Neural Network-Based Prediction Intervals and New Advances. *IEEE Transactions on Neural Networks* 22(9): 1341–1356.

Kindig, B. 2020. 5 Soon-to-Be Trends in Artificial Intelligence And Deep Learning. <https://www.forbes.com/sites/bethkindig/2020/01/31/5-soon-to-be-trends-in-artificial-intelligence-and-deep-learning>.

Knight, W. 2018. Google and Others Are Building AI Systems That Doubt Themselves AI will make better decisions by embracing uncertainty. <https://www.technologyreview.com/2018/01/09/146337/google-and-others-are-building-ai-systems-that-doubt-themselves>.



- 597 Koenker, R. W.; and Bassett, G. 1978. Regression Quantiles. 650  
598 *Econometrica* 46(1): 33–50. 651
- 599 Kwok, S. W.; and Carter, C. 1990. Multiple decision trees. In 652  
600 *Uncertainty in Artificial Intelligence*, volume 9 of *Machine* 653  
601 *Intelligence and Pattern Recognition*, 327 – 335. North-  
602 Holland.
- 603 Lakshminarayanan, B.; Pritzel, A.; and Blundell, C. 2017.  
604 Simple and Scalable Predictive Uncertainty Estimation using  
605 Deep Ensembles. In *Advances in Neural Information*  
606 *Processing Systems*, volume 30, 6402–6413. Curran Associ-  
607 ates, Inc.
- 608 Lei, J.; G’Sell, M.; Rinaldo, A.; Tibshirani, R. J.; and  
609 Wasserman, L. 2018. Distribution-Free Predictive Inference  
610 for Regression. *Journal of the American Statistical Associ-*  
611 *ation* 113(523): 1094–1111.
- 612 Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. 2018.  
613 Learning to Generalize: Meta-Learning for Domain Gener-  
614 alization. *AAAI Conference on Artificial Intelligence* .
- 615 Malinin, A.; and Gales, M. 2018. Predictive Uncertainty  
616 Estimation via Prior Networks. In *Advances in Neural In-*  
617 *formation Processing Systems*, volume 31, 7047–7058.
- 618 Neal, R. M.; et al. 2011. MCMC using Hamiltonian dynam-  
619 ics. *Handbook of markov chain monte carlo* 2(11): 2.
- 620 Nguyen, A.; Yosinski, J.; and Clune, J. 2015. Deep neural  
621 networks are easily fooled: High confidence predictions for  
622 unrecognizable images. In *IEEE Conference on Computer*  
623 *Vision and Pattern Recognition (CVPR)*, 427–436. ISSN  
624 1063-6919.
- 625 Nix, D. A.; and Weigend, A. S. 1994. Estimating the mean  
626 and variance of the target probability distribution. In *Pro-*  
627 *ceedings of 1994 IEEE International Conference on Neural*  
628 *Networks (ICNN’94)*, volume 1, 55–60. IEEE.
- 629 Osband, I.; Aslanides, J.; and Cassirer, A. 2018. Random-  
630 ized Prior Functions for Deep Reinforcement Learning. In  
631 *Advances in Neural Information Processing Systems*, vol-  
632 ume 31.
- 633 Papadopoulos, H. 2008. *Inductive Conformal Prediction:*  
634 *Theory and Application to Neural Networks*, 315–330. Cite-  
635 seer. ISBN 978-953-7619-03-9.
- 636 Redyuk, S.; Schelter, S.; Rukat, T.; Markl, V.; and Biess-  
637 mann, F. 2019. Learning to Validate the Predictions of Black  
638 Box Machine Learning Models on Unseen Data. In *Pro-*  
639 *ceedings of the Workshop on Human-In-the-Loop Data An-*  
640 *alytics*, 1–4.
- 641 Schat, E.; van de Schoot, R.; Kouw, W. M.; Veen, D.;  
642 and Mendrik, A. M. 2020. The Data Representativeness  
643 Criterion: Predicting the Performance of Supervised Clas-  
644 sification Based on Data Set Similarity. *arXiv preprint*  
645 *arXiv:2002.12105* .
- 646 Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and  
647 Salakhutdinov, R. 2014. Dropout: A Simple Way to Pre-  
648 vent Neural Networks from Overfitting. *Journal of Machine*  
649 *Learning Research* 15: 1929–1958.
- Su, J.; Vargas, D. V.; and Sakurai, K. 2019. One Pixel At- 650  
651 tack for Fooling Deep Neural Networks. *IEEE Transactions*  
652 *on Evolutionary Computation* 23(5): 828–841. ISSN 1941-  
653 0026.
- Talagala, T. S.; Li, F.; and Kang, Y. 2019. FFORMPP:  
654 Feature-based forecast model performance prediction. *arXiv*  
655 *preprint arXiv:1908.11500* . 656
- Vinayak, R. K.; and Gilad-Bachrach, R. 2015. DART:  
657 Dropouts meet Multiple Additive Regression Trees. In *Pro-*  
658 *ceedings of the Eighteenth International Conference on Ar-*  
659 *tificial Intelligence and Statistics*, volume 38 of *Proceedings*  
660 *of Machine Learning Research*, 489–497. PMLR. 661
- Vovk, V. 2012. Conditional validity of inductive conformal  
662 predictors. In *Asian conference on machine learning*, 475–  
663 490. 664
- Vovk, V.; Gammerman, A.; and Shafer, G. 2005. *Algorith-*  
665 *mic learning in a random world*. Springer Science & Busi-  
666 ness Media. 667
- Vovk, V.; Nouretdinov, I.; Manokhin, V.; and Gammerman,  
668 A. 2018. Cross-conformal predictive distributions. In *Con-*  
669 *formal and Probabilistic Prediction and Applications*, 37–  
670 51. 671
- Zadrozny, B.; and Elkan, C. 2001. Obtaining calibrated  
672 probability estimates from decision trees and naive Bayesian  
673 classifiers. In *Proceedings of the Eighteenth International*  
674 *Conference on Machine Learning*, 609–616. Morgan Kauf-  
675 mann. 676