

---

# Active Assessment of Prediction Services as Accuracy Surface Over Attribute Combinations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Our goal is to evaluate the accuracy of a black-box classification model, not as  
2 a single aggregate on a given test data distribution, but as a surface over a large  
3 number of combinations of attributes characterizing multiple test data distributions.  
4 Such attributed accuracy measures become important as machine learning models  
5 get deployed as a service, where the training data distribution is hidden from clients,  
6 and different clients may be interested in diverse regions of the data distribution.  
7 We present Attributed Accuracy Assay (AAA) — a Gaussian Process (GP)-based  
8 probabilistic estimator for such an accuracy surface. Each attribute combination,  
9 called an ‘arm’, is associated with a Beta density from which the service’s accuracy  
10 is sampled. We expect the GP to smooth the parameters of the Beta density over  
11 related arms to mitigate sparsity. We show that obvious application of GPs cannot  
12 address the challenge of heteroscedastic uncertainty over a huge attribute space that  
13 is sparsely and unevenly populated. In response, we present two enhancements:  
14 pooling sparse observations, and regularizing the scale parameter of the Beta  
15 densities. After introducing these innovations, we establish the effectiveness of  
16 AAA both in terms of its estimation accuracy and exploration efficiency, through  
17 extensive experiments and analysis.

## 18 1 Introduction

19 Increasing concentration of big data and computing resources has resulted in widespread adoption of  
20 machine learning as a service (MLaaS). The best-performing NLP, speech, image and video recog-  
21 nition tools are now provided as network services. MLaaS comes with few accuracy specifications  
22 or service level agreements, perhaps only leaderboard numbers from benchmarks that may not be  
23 closely related to most clients’ deployment data distributions. The client, therefore, finds it difficult  
24 to choose the best provider without extensive pilot trials [1]. Different clients may need to deploy the  
25 service on very different data distributions, with possibly widely different accuracy.

26 In such circumstances, we propose that a service provider, or a service standardization agency, publish  
27 the accuracy of the classifier, not as one or few aggregate numbers, but as a *surface* defined on a space  
28 of input instance *attributes* that capture the variability of consumer expectations. Indoor/outdoor,  
29 day/night, urban/rural may be attributes of input images for visual object recognition tasks. Speaker  
30 age, gender, ethnicity/accent may be attributes of input audio for speech recognition tasks. We call  
31 a combination of attributes in their Cartesian space an *arm* (borrowing from bandit terminology)  
32 <sup>1</sup>. The labeled instances used by the service provider may not represent or cover well the space of  
33 attributes of interest to subscribers. Labeled data may be proprietary and inaccessible to prospective  
34 consumers and standardization agencies. Whoever estimates the accuracy surface, therefore, needs to  
35 *actively* select instances from an unlabeled pool for labeling, presumably within a restricted budget,  
36 to adequately cover the attribute space.

---

<sup>1</sup>Figure 1 shows an example of diverse accuracy over arms.

37 Several recent studies have highlighted the variability in accuracy across data sub-populations [2, 3],  
 38 specifically in the context of fairness [4, 5, 6], and also proposed active estimation techniques of  
 39 sub-population accuracy [7, 8]. We solve a more general problem where the space of arms (sub-  
 40 population) defined by the Cartesian space of attributes grows combinatorially. This inevitably leads  
 41 to extreme sparsity of labeled instances for many arms. A central challenge is how to smooth the  
 42 estimate across related arms while faithfully representing the uncertainty for active exploration.

43 We present Attributed Accuracy Assay (AAA) — a practical system that estimates accuracy, together  
 44 with the uncertainty of the estimate, as a function of the attribute space. AAA uses these estimates  
 45 to drive the sampling policy for each attribute combination. Gaussian Process (GP) regression is a  
 46 natural choice to obtain smooth probabilistic accuracy estimates over arm attributes. However, a  
 47 straightforward GP model fails to address the challenge of heteroscedasticity that we face with uneven  
 48 and sparse supervision across arms. We model arm-specific service accuracy as drawn from a Beta  
 49 density that is characterized by mean and scale parameters, which are sampled from two GPs that are  
 50 informed by suitable trained kernels over the attribute space. We propose two further enhancements  
 51 to the training of this model. First, we recognize an over-smoothing problem with GP’s estimation  
 52 of the Beta scale parameters, and propose a Dirichlet likelihood to supervise the relative values of  
 53 scale across arms. Second, we recognize that arms with very low support interfere with learning the  
 54 kernel parameters of the GPs. We mitigate this by pooling observations across related arms. With  
 55 these fixes, AAA achieves the best estimation performance among competitive alternatives.

56 Another practical challenge in our setting is that some attributes of instances are not known exactly.  
 57 For example, attributes, such as camera shutter speed or speaker gender, may be explicitly provided  
 58 as meta information attached with instances. But other attributes, such as indoor/outdoor, or speaker  
 59 age, may have to be estimated noisily via another (attribute) classifier, because accurate human-based  
 60 acquisition of attributes would be burdensome. AAA also tackles uncertain attribute inference. Its  
 61 attribute classifiers are trained on a small amount of labeled data and their error rates are modeled in  
 62 a probabilistic framework.

63 We report on extensive experiments using several real data sets. Comparison with several estimators  
 64 based on Bernoulli arm parameters, Beta densities per arm, and even simpler forms of GPs on the  
 65 arm Beta distributions, shows that AAA is superior at quickly cutting down arm accuracy uncertainty.

66 Summarizing, our contributions are:

- 67 • We motivate and define the problem of accuracy surface estimation over a large space of attribute  
 68 combinations.
- 69 • Our proposed estimator AAA fits a Beta density for every attribute combination (arm), with its  
 70 parameters smoothed via two GPs to capture heteroscedastic uncertainty of each arm’s accuracy  
 71 under limited data settings.
- 72 • We propose two important components included in AAA: 1) a Dirichlet regularization to control  
 73 over-smoothing of the Beta scale parameters, and 2) pooled observations to reduce over-fitting of a  
 74 GP-associated kernel to sparse arms.
- 75 • We show significant gains in terms of both estimation quality and the efficiency of exploration  
 76 on four real classification models compared to existing methods. AAA obtains an average 80%  
 77 reduction in macro averaged square error over the existing methods.

## 78 2 Problem Setup

79 Our goal is to evaluate a given machine learning service model  $S$  used by a diverse set of consumers.  
 80 The service  $S : \mathcal{X} \mapsto \mathcal{Y}$  could be any predictive model that, for an input instance  $\mathbf{x} \in \mathcal{X}$ , assigns  
 81 an output label  $\hat{y} \in \mathcal{Y}$ , where  $\mathcal{Y}$  is a discrete label space. Let  $y(\mathbf{x})$  denote the true label of  $\mathbf{x}$  and  
 82  $\text{Agree}(y, \hat{y})$  denote the match between the two labels. For scalar classification,  $\text{Agree}(y, \hat{y})$  is in  
 83  $\{0, 1\}$ . For structured outputs, e.g., sequences, we could use measures like BLEU scores in  $[0, 1]$ .  
 84 Classifiers are routinely evaluated on their expected accuracy on a data distribution  $P(\mathcal{X}, \mathcal{Y})$ :

$$\rho = \mathbb{E}_{P(\mathbf{x}, y)}[\text{Agree}(y, S(\mathbf{x}))] \quad (1)$$

85 We propose to go beyond this single measure and define accuracy as a surface over a space of  
 86 attributes of the input instances. Let  $A$  denote a list of  $K$  attributes that capture the variability of  
 87 consumer expectation on which the service  $S$  will be deployed. For instance, visual object recognition  
 88 is affected by the background scene, and facial recognition is affected by demographic attributes. We  
 89 use  $A(\mathbf{x}) \in \mathcal{A}$  to denote the vector of values of attributes of input  $\mathbf{x}$  and  $\mathcal{A}$  to denote the Cartesian

90 product of the domains of all attributes. An attribute could be discrete, e.g., the ethnicity of a speaker;  
 91 Boolean, e.g., whether a scene is outdoors/indoors; or continuous, e.g., the age of the speaker in  
 92 speech recognition. Some of the attributes of  $\mathbf{x}$ , for example the camera settings of an image, may be  
 93 known exactly, and others may only be available as a distribution  $M_k(a_k|\mathbf{x})$  for an attribute  $a_k \in A$ ,  
 94 obtained from a pre-trained probabilistic classifier.

95 Generalizing from a single global expected accuracy (1), we define the accuracy surface  $\rho : \mathcal{A} \rightarrow [0, 1]$   
 96 of a service  $S$  at each attribute combination  $\mathbf{a} \in \mathcal{A}$ , given a data distribution  $P(\mathcal{X}, \mathcal{Y})$ , as

$$\rho(\mathbf{a}) = \mathbb{E}_{P(\mathbf{x}, y|A(\mathbf{x})=\mathbf{a})}[\text{Agree}(y, S(\mathbf{x}))] \quad (2)$$

97 Our goal is to provide an estimate of  $\rho(\mathbf{a})$  given two kind of data sampled from  $P(\mathcal{X}, \mathcal{Y})$ : a small  
 98 labeled sample  $D$ , and a large unlabeled sample  $U$ . In addition, we are given a budget of  $B$  instances  
 99 for which we can seek labels  $y$  from a human by selecting them from  $U$ . Applying  $M_k$  to all of  $U$  is,  
 100 however, free of cost.

101 We aim to design a probabilistic estimator for  $\rho(\mathbf{a})$ , which we denote as  $P(\rho|\mathbf{a})$  where  $\rho \in [0, 1]$   
 102 and  $\mathbf{a} \in \mathcal{A}$ . This is distinct from active learning, which selects instances to train the learner toward  
 103 greater accuracy, and also active accuracy estimation [7], which does not involve a surface over  $\mathbf{a}$ .  
 104 We also show that standard tools to regress from  $\mathbf{a}$  to  $\rho$  are worse than our proposal.

105 We measure the quality of our estimate as the square error between the gold accuracy  $\rho(\mathbf{a})$  and the  
 106 mean of the estimated accuracy distribution  $P(\rho|\mathbf{a})$ . Our estimator distribution naturally gives an  
 107 idea of the posterior variance of accuracy estimate of each attribute combination, which we use for  
 108 uncertainty-based exploration.

### 109 3 Proposed Estimator

110 We will first review recent work that leads to candidate solutions to our problem, discuss their  
 111 limitations, and finally present our solution. Initially, to keep the treatment simple, we assume  $A(\mathbf{x})$   
 112 and gold  $y$  (hence  $c = \text{Agree}(S(x), y)$ , the service correctness bit) is known for all instances. Later  
 113 in this section, we remove these assumptions.

114 The simplest option is to ignore any relationship between arms, and, for each arm  $\mathbf{a}$ , fit a suitable  
 115 density over  $\rho(\mathbf{a})$ . When this density is sampled, we get a number in  $[0, 1]$ , which is like a coin head  
 116 probability used to sample correctness bits  $c$ . For representing uncertainty of accuracy values (which  
 117 are ratios between two counts), the **Beta distribution**  $\mathfrak{B}(\cdot, \cdot)$  is a natural choice. We call this baseline  
 118 method **Beta-I**.

119 The variance of the estimated Beta density can be used for actively sampling arms. Ji et al. [7]  
 120 describe a related scenario, stressing on active sampling. However, this approach cannot share  
 121 observations or smooth the estimated density at a sparsely-populated arm with information from  
 122 similar arms. In our real-life scenario, we expect accuracy surface smoother and the number of arms  
 123 to be large enough that many arms will get very few, if any, instances.

124 The second baseline method, which we call **BernGP**, is to view the  $(\mathbf{a}, c)$  instances in  $D$  as a standard  
 125 classification data set with the binary  $c$  values as class label and  $\mathbf{a}$  as input features. Given the limited  
 126 data, we can use the well-known GP classification approach [9] for fitting smooth values  $\rho$  as a  
 127 function of  $\mathbf{a}$ . Suppose the arms  $\mathbf{a}$  can be embedded to  $\mathcal{V}(\mathbf{a})$  in a suitable space induced by some  
 128 similarity kernel. In this embedding space, we expect the accuracy of  $S$  to vary smoothly. Given a  
 129 kernel  $K_1(\mathbf{a}, \mathbf{a}')$  to guide the extent of sharing of information across arms, a standard form of this  
 130 GP would be

$$P(c|\mathbf{a}) = \text{Bernoulli}(c; \text{sigmoid}(f_{\mathbf{a}})); \quad f \sim GP(0, K_1). \quad (3)$$

131 The GP can give estimates of uncertainty of  $\rho(\mathbf{a})$ , which may be used for active sampling of arms.

132 As we will demonstrate, such GP-imposed estimate of uncertainty of  $\rho(\mathbf{a})$  is inadequate, because  
 133 it loses sight of the number of supporting observations at each arm, which could be very diverse.  
 134 This is because the standard GP assumption of homoscedasticity, that is, identical noise around each  
 135 arm is violated when observations per arm differ significantly. We therefore need a mechanism to  
 136 separately account for the uncertainty at each arm, even the unexplored ones, to guide the strategy for  
 137 actively collecting more labeled data.

138 **3.1 The basic BetaGP proposal**

139 We model arm-specific noise by allowing each arm to represent the uncertainty of  $\rho_a$ , not just by  
 140 an underlying GP as in BernGP above, but also by a separate scale parameter. Further, the scale  
 141 parameter is smoothed over neighboring arms using another GP. The influence of this scale on the  
 142 uncertainty of  $\rho_a$  is expressed by a Beta distribution as follows:

$$P(\rho|\mathbf{a}) \sim \mathfrak{B}(\rho; \phi(f_{\mathbf{a}}), \psi(g_{\mathbf{a}})) \quad (4)$$

$$\phi(f_{\mathbf{a}}) = \text{sigmoid}(f_{\mathbf{a}}), \quad f \sim GP(0, K_1), \quad (5)$$

$$\psi(g_{\mathbf{a}}) = \log(1 + e^{g_{\mathbf{a}}}), \quad g \sim GP(0, K_2), \quad (6)$$

143 where we use  $\phi(\bullet), \psi(\bullet)$  to denote the parameters of the Beta distribution at arm  $\mathbf{a}$ . The Beta  
 144 distribution is commonly represented via  $\alpha, \beta$  parameters whereas we chose the less popular mean  
 145 ( $\phi$ ) and scale ( $\psi$ ) parameters. While these two forms are functionally equivalent with  $\phi = \frac{\alpha}{\alpha+\beta}, \psi =$   
 146  $\alpha + \beta$ , we preferred the second form because imposing GP smoothness across arms on the mean  
 147 accuracy and scale seemed more meaningful. We validate this empirically in the Appendix.

148 Two kernel functions  $K_1(\mathbf{a}, \mathbf{a}'), K_2(\mathbf{a}, \mathbf{a}')$  defined over pairs of arms  $\mathbf{a}, \mathbf{a}' \in \mathcal{A}$  control the degree  
 149 of smoothness among the Beta parameters across the arms. We use an RBF kernel defined over  
 150 learned shared embeddings  $\mathcal{V}(\mathbf{a})$ :

$$K_1(\mathbf{a}, \mathbf{a}') = s_1 \exp \left[ -\frac{\|\mathcal{V}(\mathbf{a}) - \mathcal{V}(\mathbf{a}')\|^2}{l_1} \right], \quad K_2(\mathbf{a}, \mathbf{a}') = s_2 \exp \left[ -\frac{\|\mathcal{V}(\mathbf{a}) - \mathcal{V}(\mathbf{a}')\|^2}{l_2} \right] \quad (7)$$

151 where  $s_1, s_2, l_1, l_2$  denote the scale and length parameters of the two kernels. The scale and length  
 152 parameters are learned along with the parameters of embeddings  $\mathcal{V}(\mathbf{a})$  during training.

153 Initially, we assume we are given a labeled dataset  $D = \{(\mathbf{x}_i, \mathbf{a}_i, y_i) : i = 1 \dots, I\}$  with attribute  
 154 information available. Using predictions from the classification service  $S$ , we associate a 0/1 accuracy  
 155  $c_i = \text{Agree}(y_i, S(\mathbf{x}_i))$ . We can thus extend  $D$  to  $\{(\mathbf{x}_i, \mathbf{a}_i, y_i, c_i) : i \in [I]\}$ .

156 Let  $c_{\mathbf{a}} = \sum_{i:A(\mathbf{x}_i)=\mathbf{a}} c_i$  denote the total accuracy score in arm  $\mathbf{a}$ . Let  $n_{\mathbf{a}}$  denote the total number of  
 157 labeled examples in arm  $\mathbf{a}$ . The likelihood of all observations given functions  $f, g$  decomposes as a  
 158 product of Beta-binomial<sup>2</sup> distributions at each arm as follows:

$$\Pr(D|f, g) = \prod_{\mathbf{a}} \int_{\rho} \rho^{c_{\mathbf{a}}} (1 - \rho)^{n_{\mathbf{a}} - c_{\mathbf{a}}} \mathfrak{B}(\rho | \phi(f_{\mathbf{a}}), \psi(g_{\mathbf{a}})) d\rho. \quad (8)$$

$$= \prod_{\mathbf{a}} \frac{B(\phi(f_{\mathbf{a}})\psi(g_{\mathbf{a}}) + c_{\mathbf{a}}, (1 - \phi(f_{\mathbf{a}}))\psi(g_{\mathbf{a}}) + n_{\mathbf{a}} - c_{\mathbf{a}})}{B(\phi(f_{\mathbf{a}})\psi(g_{\mathbf{a}}), (1 - \phi(f_{\mathbf{a}}))\psi(g_{\mathbf{a}}))}, \quad (9)$$

159 where  $B$  is the Beta function, and the second expression is a rewrite of the **Beta-binomial likelihood**.

160 During training we calculate the posterior distribution of functions  $f, g$  using the above data likelihood  
 161  $\Pr(D|f, g)$  and GP priors given in eqns. (5) and (6). The posterior cannot be computed analytically  
 162 given our likelihood, so we use variational methods. Further, we reduce the  $\mathcal{O}(|\mathcal{A}|^3)$  complexity of  
 163 posterior computation, using the inducing point method of Hensman et al. [9], whereby we learn  
 164  $m$  locations  $\mathbf{u} \in \mathbb{R}^{d \times m}$ , mean  $\mu \in \mathbb{R}^m$ , and covariance  $\Sigma \in \mathbb{R}^{m \times m}$  of inducing points. Doing  
 165 so brings down the complexity to  $\mathcal{O}(m^2|\mathcal{A}|)$ ,  $m \ll |\mathcal{A}|$ . These parameters are learned end to end  
 166 with the parameters of the neural network used to extract embeddings  $\mathcal{V}(\mathbf{a})$  of arms  $\mathbf{a}$ , and kernel  
 167 parameters  $s_1, s_2, l_1, l_2$ . We used off-the-shelf Gaussian process library: GPyTorch [10] to train the  
 168 above likelihood with variational methods. Details of this procedure can be found in the Appendix.  
 169 We denote the posterior functions as  $P(f|D), P(g|D)$ . Thereafter, the mean estimated accuracy for  
 170 an arm  $\mathbf{a}$  is computed as

$$\mathbb{E}(\rho|\mathbf{a}) = \mathbb{E}_{f \sim P(f|D)}[\phi(f_{\mathbf{a}})]. \quad (10)$$

171 We call this setup **BetaGP**. Next, we will argue why BetaGP still has serious limitations, and offer  
 172 mitigation measures.

173 **3.2 Supervision for scale parameters**

174 We had introduced the second GP  $g_{\mathbf{a}}$  to model arm-specific noise, and similar techniques have been  
 175 proposed earlier by Lázaro-Gredilla and Titsias [11], Kersting et al. [12], Goldberg et al. [13], but  
 176 for heteroscedasticity in Gaussian observations. However, we found the posterior distribution of

<sup>2</sup>The  $\binom{n_{\mathbf{a}}}{c_{\mathbf{a}}}$  term does not apply since we are given not just counts but accuracy  $c_i$  of individual points.

177 scale values  $\psi(g_{\mathbf{a}})$  at each arm tended to converge to similar values, even across arms with orders of  
 178 magnitude difference in number of observations  $n_{\mathbf{a}}$ . On hindsight, that was to be expected, because  
 179 the data likelihood (8) increases monotonically with scale  $\psi_{\mathbf{a}}$ . The only control over its converging  
 180 to  $\infty$  is the GP prior  $g \sim GP(0, K_2)$ . In the Appendix, we illustrate this phenomenon with an  
 181 example. We propose a simple fix to the scale supervision problem. We expect the relative values of  
 182 scale across arms to reflect the distribution of the proportion of observations  $\frac{n_{\mathbf{a}}}{n}$  across arms (with  
 183  $n = \sum_{\mathbf{a}} n_{\mathbf{a}}$ ). We impose a joint Dirichlet distribution using the scale of arms  $\psi(g_{\mathbf{a}})$  as parameters,  
 184 and write the likelihood of the observed proportions as (with  $\Gamma$  denoting [Gamma function](#)):

$$\log \Pr(\{n_{\mathbf{a}}\}|g) = \sum_{\mathbf{a}} ((\psi(g_{\mathbf{a}}) - 1) \log \frac{n_{\mathbf{a}}}{n} - \log \Gamma(\psi(g_{\mathbf{a}})) + \log \Gamma(\sum_{\mathbf{a}} \psi(g_{\mathbf{a}}))) \quad (11)$$

185 We call this **BetaGP-SL**. With this as an additional term in the data likelihood, we obtained signifi-  
 186 cantly improved uncertainty estimates at each arm, as we will show in the experiment section.

### 187 3.3 Pooling for sparse observations

188 Recall that the observations are accumulation of 1/0 agreement scores for all instances that belong to  
 189 an arm. Given the nature of our problem, arms have varying levels of supervision, and also highly  
 190 varying true accuracy values. Even when the available labeled data is large, many arms will continue  
 191 to have sparse supervision because they represent rare attribute combinations. The combination  
 192 of high variance observations and sparse supervision could lead to learning of non-smooth kernel  
 193 parameters. The situation is further aggravated when learning a deep kernel. This problem has  
 194 resemblance to “collapsing variance problem” [14] such as when Gaussian Mixture models overfit on  
 195 outliers or when topic models overfit a noisy document in the corpus. Instead of depending purely  
 196 on GP priors to smooth over these noisy observations, we found it helpful to also externally smooth  
 197 noisy observations. For each arm  $\mathbf{a}$  with observations below a threshold, we mean-pool observations  
 198 from some number of nearest neighbors, weighted by their kernel similarity with  $\mathbf{a}$ . We will see that  
 199 such external smoothing resulted in significantly more accurate estimates particularly for arms with  
 200 extreme accuracy values. We call this method **BetaGP-SLP** (note that this also includes the scale  
 201 supervision objective described in the previous section). Two other mechanisms take us to the full  
 202 form of the **AAA** system, which we describe next.

### 203 3.4 Exploration

204 The variance estimate of an arm informs its uncertainty and is commonly used for efficient explo-  
 205 ration [15]. Let  $P(f|D), P(g|D)$  denote the learned posterior distribution of the GPs. Using these,  
 206 the estimated variance at an arm is given as:

$$\mathbb{V}(\rho|\mathbf{a}) = \mathbb{E}_{f \sim P(f|D), g \sim P(g|D)} \left[ \int_{\rho} (\rho - \mathbb{E}(\rho|\mathbf{a}))^2 \mathfrak{B}(\rho; \phi(f_{\mathbf{a}}), \psi(g_{\mathbf{a}})) d\rho \right] \quad (12)$$

207 where the expected value is given in eqn. (10). We use sampling to estimate the above expectation.  
 208 The arm to be sampled next is chosen as the one with the highest variance among unexplored arms.  
 209 We then sample an unexplored example with highest affiliation ( $P(\mathbf{a} | \mathbf{x})$ ) with the chosen arm.

### 210 3.5 Modeling Attribute Uncertainty

211 Recall that attributes of an instance  $\mathbf{x}$  are obtained from models  $M_k(a_k|\mathbf{x})$ ,  $k \in [K]$ , which may  
 212 be highly noisy for some attributes. Thus, we cannot assume a fixed attribute vector  $A(\mathbf{x})$  for an  
 213 instance  $\mathbf{x}$ . We address this by designing a model that can combine these noisy estimates into a  
 214 joint distribution  $P(\mathbf{a}|\mathbf{x})$  using which, we can fractionally assign each instance  $\mathbf{x}_i$  across arms. A  
 215 baseline model for  $P(\mathbf{a}|\mathbf{x})$  would be just the product  $\prod_{k=1}^K M_k(a_k|\mathbf{x})$ . However, we expect values of  
 216 attributes to be correlated (e.g. attribute ‘high-pitch’ is likely to be correlated with gender ‘female’).  
 217 Also, the probabilities  $M_k(a_k|\mathbf{x})$  may not be well-calibrated.

218 We therefore propose an alternative joint model that can both recalibrate individual classifiers via  
 219 temperature scaling [16], and model their correlation. We have a small seed labeled dataset  $D$  with  
 220 gold attribute labels, independent noisy distributions from each attribute model  $M_k(a_k|\mathbf{x})$ , and an  
 221 unlabeled dataset  $U$ . We prefer simple factorized models. We factorize  $\log \Pr(\mathbf{a}|\mathbf{x})$  as a sum of  
 222 temperature-weighted logits and a joint (log) potential as shown in expression (13) below.

$$\log \Pr(\mathbf{a}|\mathbf{x}) = \log \Pr(a_1, a_2, \dots, a_K|\mathbf{x}) = \sum_{k=1}^K t_k \log M_k(a_k|\mathbf{x}) + N(a_1, a_2, \dots, a_K) \quad (13)$$

223 Here  $N$  denotes a dense network to model the correlation between attributes, and  $t_1, \dots, t_K$  denote  
 224 the temperature parameters used to rescale noisy attribute probabilities. The maximum likelihood  
 225 over  $D$  is  $\max_{t,N} \sum_{(\mathbf{x}_i, \mathbf{a}_i) \in D} \log \Pr(\mathbf{a}_i|\mathbf{x}_i)$

$$= \max_{t,N} \sum_{\mathbf{x}_i \in D} \left\{ \sum_{k=1}^K t_k \log M_k(a_{ik}|\mathbf{x}_i) + N(a_{i1}, \dots, a_{iK}) - \log(Z_i) \right\} \quad (14)$$

226  $Z_i$  denotes the partition function for an example  $\mathbf{x}_i$  which requires summation over  $\mathcal{A}$ . It could  
 227 be intractable to compute  $Z_i$  exactly when  $\mathcal{A}$  is large. In such cases,  $Z_i$  can be approximated by  
 228 sampling. In our case, we could get exact estimates.

229 In addition to  $D$ , we use the unlabeled instances  $U$  with predictions from attribute predictors filling  
 230 the role of gold-attributes. Details on how we train the parameters on large but noisy  $U$  and small but  
 231 correct  $D$  can be found in the Appendix.

232 The estimation method of BetaGP-SLP with variance based exploration and calibration described here  
 233 constitute our proposed estimator: AAA. Detailed pseudo-code of AAA is given in the Appendix.

## 234 4 Experiments

235 Our exploration of various methods and data sets is guided by the following research questions.

- 236 • How do various methods for arm accuracy estimation compare?
- 237 • To what extent do BetaGP, scale supervision and pooled observations help beyond BernGP?
- 238 • For the best techniques from above, how do various active exploration strategies compare?
- 239 • How well does our proposed model of attribute uncertainty work?

### 240 4.1 Data sets and tasks

241 We experiment with two real data sets and tasks. Our two tasks are male-female gender classification  
 242 with two classes and animal classification with 10 classes.

243 **Male-Female classification (MF):** CelebA [17] is a popular celebrity faces and attribute data set  
 244 which identifies the gender of celebrities among 39 other binary attributes. The label is gender. The  
 245 accuracy surface spans various demographic, style, and personality related attributes. We hand-pick  
 246 a subset of 12 attributes that we deem important for gender classification. Gender-neural attributes  
 247 such as wearing spectacles or hat are ignored (see Appendix for more details). A subset of 50,000  
 248 examples is used to train classifiers on each of the 12 attributes using a pretrained ResNet-50 model.  
 249 The remaining 150,000 examples in the data set are set as the unlabeled pool from which we actively  
 250 explore new examples for human feedback.

251 **Animal classification (AC):** COCO-Stuff [18] provides an image collection. For each image,  
 252 labels for foreground (cow, camel) and background (sky, snow, water) ‘stuff’ are available. Visual  
 253 recognition models often correlate the background scene with the animal label such as camel  
 254 with deserts and cow with meadows. Thus, foreground stuff labels are our regular  $y$ -labels while  
 255 background stuff labels supply our notion of attributes.

256 We collapse fine stuff labels into five coarse labels using the dataset provided label hierarchy. These  
 257 are: water, ground, sky, structure, furniture. The Coco dataset has around 90 object labels. Here  
 258 we use a subset of 10 labels corresponding to animals. We take special care to filter out images  
 259 with multiple/no animals and adapt the pixel segmentation/classification task to object classification  
 260 (see the Appendix for more details). The image is further annotated with the five binary labels  
 261 corresponding to five coarse stuff labels. The scene descriptive five binary labels and ten object labels  
 262 make up for  $32 \times 10 = 320$  attribute combinations.

### 263 4.2 Service Models

264 For the MF task, we use two service models  $S$ . **MF-CelebA** is a service model for gender classifica-  
 265 tion. To simulate separate  $D$  and  $U$ , it is trained on a random subset of CelebA with a ResNet-50

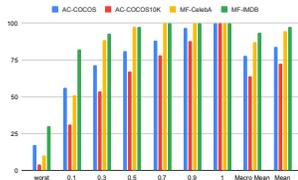


Figure 1: Macro and micro averaged accuracy (right most) and ten quantiles (x-axis) of per-arm accuracy (y-axis).

Service→	AC-COCOS10K	AC-COCOS	MF-IMDB	MF-CelebA
CPredictor	5.4 / 15.0	3.2 / 9.4	1.2 / 8.2	5.2 / 35.9
Beta-I	7.0 / 15.6	4.3 / 10.0	1.6 / 8.4	4.7 / 30.3
BernGP	7.0 / 13.2	3.5 / 8.6	1.7 / 7.6	4.9 / 28.1
BetaGP	7.1 / 14.3	3.3 / 7.9	2.2 / 6.6	4.6 / 25.9
BetaGP-SL	5.3 / 11.7	2.8 / 6.8	1.4 / 4.4	4.1 / 22.6
BetaGP-SLP	4.7 / 10.4	2.8 / 5.7	1.4 / 3.9	4.3 / 23.3

Table 2: Comparing different estimation methods on labeled data size 2000 across four tasks. No exploration is involved. Each cell shows two numbers in the format “macro MSE / worst MSE” obtained over three runs. BetaGP-SLP generally gives the lowest MSE.

266 model. **MF-IMDB** is a publicly available<sup>3</sup> classifier trained on IMBD-Wiki dataset, also using the  
 267 ResNet50 architecture. The attribute predictors are trained using ResNet-50 on a subset of the CelebA  
 268 dataset for both service models.

269 For the AC task, we use two publicly available<sup>4</sup> service models  $S$ . **AC-COCOS** was trained on  
 270 COCOS data set with 164K examples. **AC-COCOS10k** was trained on COCOS10K, an earlier  
 271 version of COCOS with only 10K instances. We use these architectures for both label and attribute  
 272 prediction. See Appendix for more details on attribute predictor, service models and their architecture.  
 273 In Figure 1, we illustrate some statistics of the shape of the accuracy surface for the four dataset-task  
 274 combinations. Although  $S$ ’s mean accuracy (right most bars) is reasonably high, the accuracy of the  
 275 arms in the 10% quantile is abysmally low, while arms in the top quantiles have near perfect accuracy.  
 276 This further motivates the need for an accuracy surface instead of single accuracy estimate.

### 277 4.3 Methods Compared

278 We compare the proposed estimation method AAA against natural baselines, alternatives, and  
 279 ablations. Some of the methods, such as **Beta-I**, **BernGP** and **BetaGP**, we have already defined in  
 280 Section 3. We train methods BernGP and BetaGP using the default arm-level likelihood. We also  
 281 separately evaluate the impact of our fixes on BetaGP with only scale supervision: **BetaGP-SL** and  
 282 along with mean pooling: **BetaGP-SLP**. We also include a trivial baseline: **CPredictor** which fits  
 283 all the arms with a global accuracy estimated using gold  $D$ . We do not try sparse observation pooling  
 284 with Beta-I since there is no notion of per-arm closeness. We also skip it on BernGP since it is worse  
 285 than BetaGP as we will show below.

### 286 4.4 Other experimental settings

287 **Gold accuracies  $\rho(a)$ :** We compute the oracular accuracy per arm using the gold attribute/label  
 288 values of examples in  $U$  which we treat as unlabeled during exploration. For every arm with at  
 289 least five examples, we set its accuracy to be the empirical estimate obtained through the average  
 290 correctness of all the examples that belong to the arm. We discard and not evaluate on any arms with  
 291 fewer than five examples since their true accuracy cannot reliably be estimated.

292 **Warm start:** We start with 500 examples having gold attributes+labels to warm start all our experi-  
 293 ments. The random seed also picks this random subset of 500 labeled examples. We calculate the  
 294 overall accuracy of the classifier on these warm start examples as  $\hat{\rho} = (\sum_i c_i) / (\sum_i 1)$ . For all arms  
 295 we use a default smoothing to  $[\lambda\hat{\rho}, \lambda]$  where  $\lambda = 0.1$ , a randomly picked low value.

296 Unless otherwise specified, we give equal importance to each arm and report MSE macroaveraged  
 297 over all arms. Along with macro MSE, we also sometimes report MSE on the subset of 50 worst  
 298 accuracy arms, referred to as worst MSE. We report other aggregate errors in the Appendix. All the  
 299 numbers reported here are averaged over three runs with different random seeds. The initial set of  
 300 warm-start examples ( $D$ ) is also changed between the runs. In the case of BetaGP-SLP, for any arm  
 301 with observation count below 5, we mean pool from its three closest neighbours.

302 In the following Sections: 4.5 and 4.6, we compare various estimation and exploration strategies  
 303 with  $P(a|x)$  noise calibrated as described in Section 3.5. In Section 4.7, we study different forms of  
 304 calibration and demonstrate the superiority of our proposed calibration technique of Equation (13).

<sup>3</sup><https://github.com/yu4u/age-gender-estimation>

<sup>4</sup><https://github.com/kazuto1011/deeplab-pytorch/>

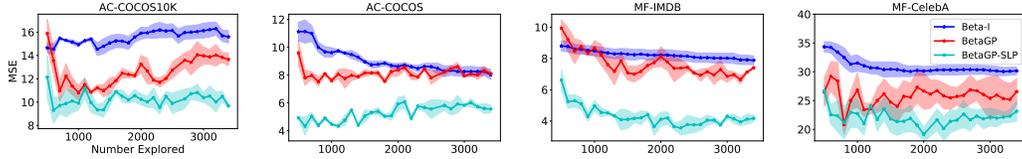


Figure 3: Comparison of estimation methods using worst MSE metric. The shaded region shows standard error. BetaGP-SLP consistently performs better than BetaGP. Beta-I is worse than its smoother counterparts.

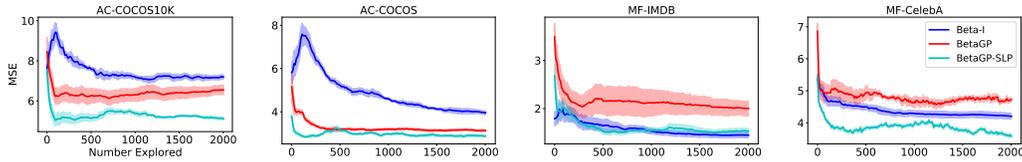


Figure 4: Comparison of exploration methods. BetaGP-SLP reduces macro MSE fastest most of the time. Shaded region shows standard error.

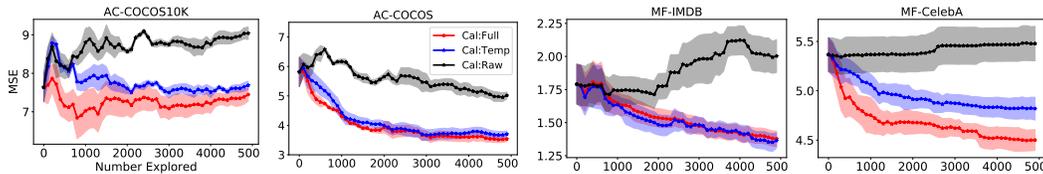


Figure 5: Calibration methods compared on different tasks. Cal:Full (red) includes temperature-based recalibration and correlation modeling with joint potential and gives the best macro MSE. Shaded region shows standard error.

#### 305 4.5 Accuracy Estimation Quality

306 We evaluate methods on their estimation quality when each method is provided with exactly the same (randomly chosen) labeled set. We compare the four service models when fitted on labeled  
 307 data of size 2,000 and the results appear in Table 2. Note that we only have label supervision on  $\mathcal{Y}$   
 308 in the labeled data. Table 2 shows macro and worst MSE, standard deviation for each metric can  
 309 be found in Appendix. In Figure 3, we show worst MSE for a range of labeled data sizes along  
 310 with their error bars. We make the following observations. **Smoothing helps:** Since we have a large  
 311 number of arms, we expect Beta-I to fare worse than its smooth counterparts (BernGP and BetaGP),  
 312 especially on the worst arms. This is confirmed in the table. In three out of four cases, this method  
 313 is worse than even the constant predictor CPredictor on both metrics. **Modeling arm specific noise**  
 314 **helps:** BetaGP is better than BernGP on almost all the cases in the table. **Significant gains when**  
 315 **the scale supervision problem of BetaGP is fixed:** BetaGP-SL is significantly better than BetaGP  
 316 in the table and figure. **Our pooling strategy helps:** BetaGP-SLP improves BetaGP-SL over worst  
 317 MSE without hurting macro MSE as seen in the table and figure.  
 318

#### 319 4.6 Exploration Efficiency

320 We compare different methods that use their own estimated variance for exploring instances to  
 321 label (Section 3.4), as a function of the number of explored examples — see Figure 4. In most  
 322 cases, BetaGP-SLP gives the smallest macro MSE, beating Beta-I and BetaGP. Note Beta-I is the  
 323 exploration method recently suggested in [7]. We observe that BetaGP provides very poor exploration  
 324 quality, indicating that the uncertainty of arms is not captured well by just using two GPs. In fact,  
 325 in many cases BetaGP is worse than Beta-I, even though we saw the opposite trend in estimation  
 326 quality (Figure 3). These experiments brings out the significant role of Dirichlet scale supervision  
 327 and pooled observations in enhancing the uncertainty estimates at each arm.

#### 328 4.7 Impact of Calibration

329 We consider two baselines along with our method explained in Section 3.5: **Cal:Raw**, which uses the  
 330 predicted attribute from the attribute models without any calibration and **Cal:Temp**, which calibrates

331 only the temperature parameters shown in eqn. (13), i.e., without the joint potential part. We refer  
332 to our method of calibration using temperature and joint potentials as **Cal:Full**. We compare these  
333 on the four tasks with estimation method set to Beta-I and random exploration strategy. Figure 5  
334 compares the three methods: Cal:Raw(Black), Cal:Temp(Blue), Cal:Full(Red). The X-axis is the  
335 number of explored examples beyond  $D$ , and Y-axis is estimation error. Observe how Cal:Temp and  
336 Cal:Full are consistently better than Cal:Raw, and Cal:Full is better than Cal:Temp.

## 337 5 Related Work

338 Our problem of actively estimating the accuracy *surface* of a classifier generalizes the more estab-  
339 lished problem of estimating a single accuracy *score* [19, 20, 21, 22, 23, 24]. For that problem, a  
340 known solution is stratified sampling, which partitions data into homogeneous strata and then seeks  
341 examples from regions with highest uncertainty and support. If we view each arm as a stratum, our  
342 method follows similar strategy. A key difference in our setting is that low support arms cannot be  
343 ignored. This makes it imperative to calibrate well the uncertainty under limited and skewed support  
344 distribution. The setting of Ji et al. [7] is the closest to ours. However, their work only considers a  
345 single attribute which they fit using Beta-I, whereas we focus on the challenges of estimating accuracy  
346 over many sparsely populated attribute combinations.

347 **Sub-population performance:** Several recent papers have focused on identifying sub-populations  
348 with significantly worse accuracy than aggregated accuracy [2, 3, 6, 8, 25, 26]. Some of these have  
349 also proposed sample-efficient techniques [6, 8] for estimation of performance on specific sub-groups,  
350 such as the ones defined by attributes like gender and race. Our accuracy surface estimation problem  
351 can be seen as a generalization where we need to estimate for all sub-groups defined in the Cartesian  
352 space of pre-specified semantic attributes. Mitchell et al. [5] recommend enclosing *model cards* with  
353 released or deployed models. In model cards, they suggest reporting performance under various  
354 relevant demographic/environmental factors which resembles the accuracy surface.

355 **Experiment design:** Another related area is experiment design using active explorations with GPs  
356 [27]. Their goal is to find the mode of the surface whereas our goal is to estimate the entire surface.  
357 Further, each arm in our setting corresponds to multiple instances, which gives rise to a degree of  
358 heteroscedasticity and input-dependent noise that is not modeled in their settings. Lázaro-Gredilla  
359 and Titsias [11], Kersting et al. [12] propose to handle heteroscedasticity by using a separate GP  
360 to model the variance at each arm. However, we showed the importance of additional terms in our  
361 likelihood and observation pooling to reliably represent estimation uncertainty. Wenger et al. [28]  
362 propose observation pooling for estimating smooth Betas but they assume a fixed kernel.

363 **Model debugging:** Testing deep neural network (DNN) is another emerging area [29]. Pei et al.  
364 [30], Tian et al. [31], Sun et al. [32], Odena et al. [33] propose to generate test examples with good  
365 coverage over all activations of a DNN. Ribeiro et al. [34], Kim et al. [35] identify rules that explain  
366 the model predictions.

## 367 6 Conclusion

368 We presented AAA, a new approach to estimate the accuracy of a classification service, not as a  
369 single number, but as a surface over a space of attributes (arms). AAA models uncertainty with a  
370 Beta distribution at each arm and regresses these parameters using two Gaussian Processes to capture  
371 smoothness and generalize to unseen arms. We proposed an additional Dirichlet likelihood to mitigate  
372 an over-smoothing problem with GP’s estimation of Beta distributions’ scale parameters. Further, to  
373 protect these high-capacity GPs from unreliable accuracy observations at sparsely populated arms,  
374 we propose to use an observation pooling strategy. Finally, we show how to handle noisy attribute  
375 labels by an efficient joint recalibration method. Evaluation on real-life datasets and classification  
376 services show the efficacy of AAA, both in estimation and exploration quality.

377 **Limitation and future work:** (1) We have evaluated AAA on the order of thousands of arms. Even  
378 larger attribute spaces could unearth more challenges. (2) Identifying relevant attributes for an  
379 application can be non-trivial. Future work could devise strategies for attribute selection. (3) It may  
380 be hard to characterize test-time data shifts, particularly for text — there could be subtle changes in  
381 word usage, style, or punctuation. A more expressive attribute space needs to be developed for text  
382 applications.

383 **References**

- 384 [1] Lingjiao Chen, Matei Zaharia, and James Zou. Frugalml: How to use ml prediction apis more  
385 accurately and cheaply. *arXiv preprint arXiv:2006.07512*, 2020.
- 386 [2] Adarsh Subbaswamy, Roy Adams, and Suchi Saria. Evaluating model robustness and stability  
387 to dataset shift. In *International Conference on Artificial Intelligence and Statistics*, pages  
388 2611–2619. PMLR, 2021.
- 389 [3] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally  
390 robust neural networks for group shifts: On the importance of regularization for worst-case  
391 generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- 392 [4] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in  
393 commercial gender classification. In *Conference on fairness, accountability and transparency*,  
394 pages 77–91. PMLR, 2018.
- 395 [5] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchin-  
396 son, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting.  
397 In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229,  
398 2019.
- 399 [6] Disi Ji, Padhraic Smyth, and Mark Steyvers. Can i trust my fairness metric? assessing fairness  
400 with unlabeled data and bayesian inference. *arXiv preprint arXiv:2010.09851*, 2020.
- 401 [7] Disi Ji, Robert L Logan IV, Padhraic Smyth, and Mark Steyvers. Active bayesian assessment  
402 for black-box classifiers. *arXiv preprint arXiv:2002.06532*, 2020. URL [https://arxiv.org/  
403 pdf/2002.06532](https://arxiv.org/pdf/2002.06532).
- 404 [8] Andrew C Miller, Leon A Gatys, Joseph Futoma, and Emily B Fox. Model-based metrics:  
405 Sample-efficient estimates of predictive model subpopulation performance. *arXiv preprint  
406 arXiv:2104.12231*, 2021.
- 407 [9] James Hensman, Alexander Matthews, and Zoubin Ghahramani. Scalable variational gaussian  
408 process classification. *JMLR*, 2015.
- 409 [10] Jacob R Gardner, Geoff Pleiss, David Bindel, Kilian Q Weinberger, and Andrew Gordon  
410 Wilson. Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration.  
411 In *Advances in Neural Information Processing Systems*, 2018.
- 412 [11] Miguel Lázaro-Gredilla and Michalis K. Titsias. Variational heteroscedastic gaussian process  
413 regression. In *ICML*, 2011.
- 414 [12] K. Kersting, C. Plagemann, P. Pfaff, and W. Burgard. Most likely heteroscedastic gaussian  
415 process regression. In *ICML '07*, 2007.
- 416 [13] Paul W. Goldberg, Christopher K. I. Williams, and Christopher M. Bishop. Regression with  
417 input-dependent noise: A gaussian process treatment. In *Proceedings of the 10th International  
418 Conference on Neural Information Processing Systems, NIPS'97*, 1997.
- 419 [14] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- 420 [15] Eric Schulz, Maarten Speekenbrink, and Andreas Krause. A tutorial on gaussian process  
421 regression: Modelling, exploring, and exploiting functions. *Journal of Mathematical Psychology*,  
422 85:1–16, 2018.
- 423 [16] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural  
424 networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML  
425 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 1321–1330, 2017.
- 426 [17] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the  
427 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 428 [18] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in  
429 context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
430 pages 1209–1218, 2018.
- 431 [19] Christoph Sawade, Niels Landwehr, Steffen Bickel, and Tobias Scheffer. Active risk estimation.  
432 In *ICML*, 2010.
- 433 [20] Christoph Sawade, Niels Landwehr, and Tobias Scheffer. Active comparison of prediction  
434 models. In *Advances in Neural Information Processing Systems*, volume 25, pages 1754–1762,  
435 2012.

- 436 [21] Namit Katariya, Arun Iyer, and Sunita Sarawagi. Active evaluation of classifiers on large  
437 datasets. In *ICDM*, 2012.
- 438 [22] Gregory Druck and Andrew McCallum. Toward interactive training and evaluation. In *CIKM*,  
439 2011.
- 440 [23] Paul N. Bennett and Vitor R. Carvalho. Online stratified sampling: evaluating classifiers at  
441 web-scale. In *CIKM*, 2010.
- 442 [24] Mohammad Reza Karimi, Nezihe Merve Gürel, Bojan Karlas, Johannes Rausch, Ce Zhang,  
443 and Andreas Krause. Online active model selection for pre-trained classifiers. *CoRR*,  
444 abs/2010.09818, 2020.
- 445 [25] Luke Oakden-Rayner, Jared Dunnmon, Gustavo Carneiro, and Christopher Ré. Hidden strat-  
446 ification causes clinically meaningful failures in machine learning for medical imaging. In  
447 *Proceedings of the ACM conference on health, inference, and learning*, pages 151–159, 2020.
- 448 [26] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay  
449 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds:  
450 A benchmark of in-the-wild distribution shifts. *arXiv preprint arXiv:2012.07421*, 2020.
- 451 [27] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian pro-  
452 cess optimization in the bandit setting: No regret and experimental design. *arXiv preprint*  
453 *arXiv:0912.3995*, 2009.
- 454 [28] Jonathan Wenger, Hedvig Kjellström, and Rudolph Triebel. Non-parametric calibration for  
455 classification. In *International Conference on Artificial Intelligence and Statistics*, pages  
456 178–190. PMLR, 2020.
- 457 [29] J. M. Zhang, M. Harman, L. Ma, and Y. Liu. Machine learning testing: Survey, landscapes and  
458 horizons. *IEEE Transactions on Software Engineering*, pages 1–1, 2020.
- 459 [30] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. Deepxplore: Automated whitebox  
460 testing of deep learning systems. In *proceedings of the 26th Symposium on Operating Systems*  
461 *Principles*, pages 1–18, 2017.
- 462 [31] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: Automated testing of deep-  
463 neural-network-driven autonomous cars. In *Proceedings of the 40th international conference*  
464 *on software engineering*, pages 303–314, 2018.
- 465 [32] Youcheng Sun, Min Wu, Wenjie Ruan, Xiaowei Huang, Marta Kwiatkowska, and Daniel  
466 Kroening. Concolic testing for deep neural networks. In *Proceedings of the 33rd ACM/IEEE*  
467 *International Conference on Automated Software Engineering*, pages 109–119, 2018.
- 468 [33] Augustus Odena, Catherine Olsson, David Andersen, and Ian Goodfellow. Tensorfuzz: Debug-  
469 ging neural networks with coverage-guided fuzzing. In *International Conference on Machine*  
470 *Learning*, pages 4901–4911, 2019.
- 471 [34] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-  
472 agnostic explanations. In *AAAI*, volume 18, pages 1527–1535, 2018.
- 473 [35] Edward Kim, Divya Gopinath, Corina Pasareanu, and Sanjit A Seshia. A programmatic and  
474 semantic approach to explaining and debugging neural network based object detectors. In  
475 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
476 11128–11137, 2020.

## 477 Checklist

- 478 1. For all authors...
- 479 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
480 contributions and scope? [Yes]
- 481 (b) Did you describe the limitations of your work? [Yes] In Section 6.
- 482 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We do not  
483 foresee any negative societal impact of our our work.
- 484 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?  
485 [Yes]
- 486 2. If you are including theoretical results...
- 487 (a) Did you state the full set of assumptions of all theoretical results? [N/A]

- 488 (b) Did you include complete proofs of all theoretical results? [N/A]
- 489 3. If you ran experiments...
- 490 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
- 491 results (either in the supplemental material or as a URL)? [Yes] They are included in the
- 492 supplementary material.
- 493 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
- 494 chosen)? [Yes] They have been discussed in sufficient detail in Section 4 and more details
- 495 are provided in Appendix.
- 496 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
- 497 multiple times)? [Yes] All our experiments are reported from multiple seeds. All our plots
- 498 show error bar and std. dev. for all the experiments can be found in the Appendix.
- 499 (d) Did you include the total amount of compute and the type of resources used (e.g., type of
- 500 GPUs, internal cluster, or cloud provider)? [N/A] We do not see them as relevant for our
- 501 paper, especially since the computation is cheap. We noted asymptotic complexity of the
- 502 algorithm when needed.
- 503 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 504 (a) If your work uses existing assets, did you cite the creators? [Yes] All our datasets are publicly
- 505 available and are cited, noted in the Section 4.
- 506 (b) Did you mention the license of the assets? [Yes] The pointers to the dataset contain the
- 507 license
- 508 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A] We
- 509 do not propose any new assets.
- 510 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 511 using/curating? [N/A]
- 512 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 513 information or offensive content? [N/A] All our datasets are standard and no violations of
- 514 these kind have been reported on these datasets.
- 515 5. If you used crowdsourcing or conducted research with human subjects...
- 516 (a) Did you include the full text of instructions given to participants and screenshots, if applica-
- 517 ble? [N/A]
- 518 (b) Did you describe any potential participant risks, with links to Institutional Review Board
- 519 (IRB) approvals, if applicable? [N/A]
- 520 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on
- 521 participant compensation? [N/A]