

---

# Reliable Post hoc Explanations: Modeling Uncertainty in Explainability

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 As black box explanations are increasingly being employed to establish model  
2 credibility in high stakes settings, it is important to ensure that these explanations  
3 are accurate and reliable. However, prior work demonstrates that explanations  
4 generated by state-of-the-art techniques are inconsistent, unstable, and provide  
5 very little insight into their correctness and reliability. In addition, these methods  
6 are also computationally inefficient, and require significant hyper-parameter tuning.  
7 In this paper, we address the aforementioned challenges by developing a novel  
8 Bayesian framework for generating local explanations along with their associated  
9 uncertainty. We instantiate this framework to obtain Bayesian versions of LIME and  
10 KernelSHAP which output credible intervals for the feature importances, capturing  
11 the associated uncertainty. The resulting explanations not only enable us to make  
12 concrete inferences about their quality (e.g., there is a 95% chance that the feature  
13 importance lies within the given range), but are also highly consistent and stable.  
14 We carry out a detailed theoretical analysis that leverages the aforementioned  
15 uncertainty to estimate how many perturbations to sample, and how to sample for  
16 faster convergence. This work makes the first attempt at addressing several critical  
17 issues with popular explanation methods in one shot, thereby generating consistent,  
18 stable, and reliable explanations with guarantees in a computationally efficient  
19 manner. Experimental evaluation with multiple real world datasets and user studies  
20 demonstrate that the efficacy of the proposed framework.

## 21 1 Introduction

22 As machine learning (ML) models get increasingly deployed in domain such as healthcare and  
23 criminal justice, it is important to ensure that decision makers have a clear understanding of the  
24 behavior of these models. However, ML models that achieve state-of-the-art accuracy are typically  
25 complex *black boxes* that are hard to understand. As a consequence, there has been a surge in post  
26 hoc techniques for explaining black box models [1–10]. Most popular among these techniques are  
27 local explanation methods which explain complex black box models by constructing interpretable  
28 local approximations (e.g., LIME [2], SHAP [4], MAPLE [11], Anchors [1]). Due to their generality,  
29 these methods are being leveraged to explain a number of classifiers including deep neural networks  
30 and ensemble models in a variety of domains such as law, medicine, and finance [12, 13].

31 Existing local explanation methods, however, suffer from several drawbacks. Explanations generated  
32 using these methods may be unstable [14–18], i.e., negligibly small perturbations to an instance  
33 can result in substantially different explanations. These methods are also inconsistent [19] i.e.,  
34 multiple runs on the same input instance with the same parameter settings may result in vastly  
35 different explanations. There are also no reliable metrics to ascertain the quality of the explanations  
36 output by these methods. Commonly used metrics such as explanation fidelity rely heavily on the

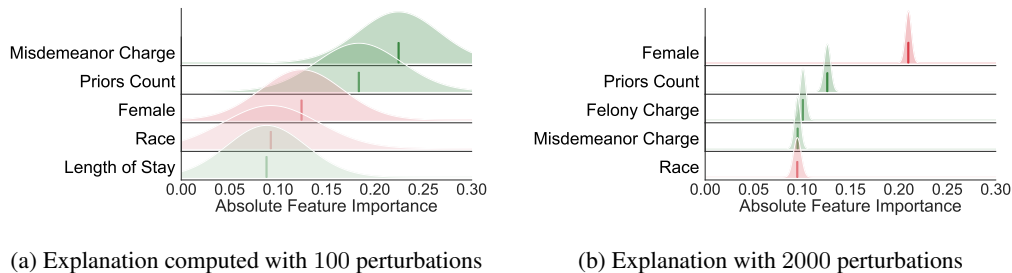


Figure 1: **Example explanations** on for an instance from the COMPAS dataset, where vertical lines indicate the feature importance by LIME (red is negative effect, green is positive) and the shaded region visualizes the uncertainty estimated by BayesLIME. While LIME produces very different and contradictory feature importance for different number of perturbations (1a and 1b), BayesLIME provides more context. The overlapping uncertainty intervals in the explanation computed with 100 perturbations (1a) indicate that it is unclear which feature is the most important. However, the tighter uncertainty intervals in the explanation computed with 2K perturbations (1b) clearly indicates that Female is the most important.

37 implementation details of the explanation method (e.g., the perturbation function used in LIME)  
 38 and do not provide a true picture of the explanation quality [20]. Furthermore, there exists little to  
 39 no guidance on determining the values of certain hyperparameters that are critical to the quality of  
 40 the resulting local explanations (e.g., number of perturbations in case of LIME). Local explanation  
 41 methods are also computationally inefficient i.e., they typically require a large number of black box  
 42 model queries to construct local approximations [21]. This can be prohibitively slow especially in  
 43 case of complex neural models.

44 In this paper, we identify that modeling uncertainty in black box explanations is the key to addressing  
 45 all the aforementioned challenges. To this end, we propose a novel Bayesian framework for generating  
 46 local explanations along with their associated uncertainty. We instantiate this framework to obtain  
 47 Bayesian versions of LIME and KernelSHAP, namely BayesLIME and BayesSHAP, that not only  
 48 output point-wise estimates of feature importance but also their associated uncertainty in the form  
 49 of credible intervals (See Figure 1). We derive closed form expressions for the posteriors of the  
 50 explanations thereby eliminating the need for any additional computational complexity. The credible  
 51 intervals produced by our framework not only allow us to make concrete inferences about the  
 52 quality of the resulting explanations but also produce explanations that satisfy user specified levels  
 53 of uncertainty (e.g., an end user may request for explanations that satisfy a certain 95% confidence  
 54 level). In addition, the resulting explanations are also highly consistent and stable. *To the best of*  
 55 *our knowledge, this work makes the first attempt at addressing several critical challenges in popular*  
 56 *explanation methods in one-shots, thereby generating consistent, stable, and reliable explanations*  
 57 *with guarantees in a computationally efficient manner.*

58 We carry out theoretical analysis that leverages the measures of uncertainty (credible intervals)  
 59 produced by our framework to estimate the values of critical hyperparameters. More specifically, we  
 60 derive a closed form expression for the number of perturbations required to generate explanations  
 61 that satisfy desired levels of confidence. We also propose a novel sampling technique called *focused*  
 62 *sampling* that leverages uncertainty to determine how to sample perturbations for faster convergence,  
 63 thereby enabling our framework to generate explanations in a computationally efficient manner.

64 We evaluate the efficacy of the proposed framework on a variety of datasets including COMPAS,  
 65 German Credit, ImageNet, and MNIST. Our results demonstrate that the explanations output by our  
 66 framework are not only highly reliable, but also very consistent and stable (53% more stable than  
 67 LIME/SHAP on an average). Our experimental results also confirm that we can accurately estimate  
 68 the number of perturbations needed to generate explanations with a desired level of uncertainty, and  
 69 that our uncertainty sampling technique speeds up the process of generating explanations by up to a  
 70 factor of 2 relative to random sampling of perturbations. Lastly, we carry out a user study with 31  
 71 human subjects to evaluate the quality of the explanations generated by our framework, demonstrating  
 72 that our explanations accurately capture the importance of the most influential features.

73 **2 Notation & Background**

74 Here we introduce notation and discuss two relevant prior approaches, LIME and KernelSHAP.

75 **Notation** Let  $f : \mathbb{R}^d \rightarrow [0, 1]$  denote a black box classifier that takes a data point  $x$  with  $d$  features,  
 76 and returns the *probability* that  $x$  belongs to a certain class. Our goal is to explain individual  
 77 predictions of  $f$ . Let  $\phi \in \mathbb{R}^d$  denote the explanation in terms of feature importances for the prediction  
 78  $f(x)$ , i.e. coefficients  $\phi$  are treated as the feature *contributions* to the black box prediction. Note  
 79 that  $\phi$  captures the coefficients of a linear model. Let  $\mathcal{Z}$  be a set of  $N$  randomly sampled instances  
 80 (perturbations) around  $x$ . The proximity between  $x$  and any  $z \in \mathcal{Z}$  is given by  $\pi_x(z) \in \mathbb{R}$ . We denote  
 81 the vector of these distances over the  $N$  perturbations in  $\mathcal{Z}$  as  $\Pi_x(\mathcal{Z}) \in \mathbb{R}^N$ . Let  $Y \in [0, 1]$  be the  
 82 vector of the black box predictions  $f(z)$  corresponding to each of the  $N$  instances in  $\mathcal{Z}$ .

83 **LIME** [2] and **KernelSHAP** [4] are popular *model-agnostic local explanation* approaches that  
 84 explain predictions of a classifier  $f$  by learning a linear model  $\phi$  locally around each prediction (i.e.  
 85  $y \sim \phi^T z$ ). The objective function for both LIME and KernelSHAP constructs an explanation that  
 86 approximates the behavior of the black box accurately in the vicinity (neighborhood) of  $x$ .

$$\arg \min_{\phi} \sum_{z \in \mathcal{Z}} [f(z) - \phi^T z]^2 \pi_x(z). \quad (1)$$

87 The above objective function has the following closed form solution:

$$\hat{\phi} = (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + \mathbb{I})^{-1} (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) Y) \quad (2)$$

88 The main difference between LIME and KernelSHAP lies in how  $\pi_x(z)$  is chosen. In LIME, it is  
 89 chosen heuristically:  $\pi_x(z)$  is computed as the cosine or  $l_2$  distance. KernelSHAP leverages game  
 90 theoretic principles to compute  $\pi_x(z)$ , guaranteeing that explanations satisfy certain properties.

91 **3 Our Framework: Bayesian Local Explanations**

92 In this section, we introduce our Bayesian framework which is designed to capture the uncertainty  
 93 associated with local explanations of black box models. First, we discuss the generative process and  
 94 inference procedure for the framework. Then, we highlight how our framework can be instantiated  
 95 to obtain Bayesian versions of LIME and SHAP. Lastly, we present detailed theoretical analysis for  
 96 estimating the values of critical hyperparameters, and discuss how to efficiently construct highly  
 97 accurate explanations with uncertainty guarantees using our framework.

98 **3.1 Constructing Bayesian Local Explanations**

99 Our goal here is to explain the behavior of a given black box model  $f$  in the vicinity of an instance  
 100  $x$  while also capturing the uncertainty associated with the explanation. To this end, we propose  
 101 a Bayesian framework for constructing local linear model based explanations and capturing their  
 102 associated uncertainty. We model the black box prediction of each perturbation  $z$  as a linear  
 103 combination of the corresponding feature values ( $\phi^T z$ ) plus an error term ( $\epsilon$ ) as shown in Eqn (3).  
 104 While the weights of the linear combination  $\phi$  capture the feature importances and thereby constitute  
 105 our explanation,  $\epsilon$  captures the error that arises due to the mismatch between our explanation  $\phi$  and  
 106 the local decision surface of the black box model  $f$ . Our complete generative process is shown below:

$$y|z, \phi, \epsilon \sim \phi^T z + \epsilon \quad \epsilon \sim \mathcal{N}(0, \frac{\sigma^2}{\pi_x(z)}) \quad \phi | \sigma^2 \sim \mathcal{N}(0, \sigma^2 \mathbb{I}) \quad \sigma^2 \sim \text{Inv-}\chi^2(n_0, \sigma_0^2). \quad (3)$$

107 The error term is modeled as a Gaussian whose variance relies on the proximity function  $\pi_x(z)$  i.e.,  
 108  $\epsilon \sim \mathcal{N}(0, \frac{\sigma^2}{\pi_x(z)})$ . This proximity function ensures that perturbations closer to the data point  $x$  are  
 109 modeled accurately, while allowing more room for error in case of perturbations that are farther away.  
 110  $\pi_x(z)$  can be computed using cosine or  $l_2$  distance or other game theoretic principles similar to that  
 111 of LIME and KernelSHAP (see Section 2). The conjugate priors on  $\phi$  and  $\sigma^2$  are shown in Eqn (3).

112 Thus, our generative process corresponds to the Bayesian version of the weighted least squares for-  
 113 mulation of LIME and KernelSHAP outlined in Eqn. (1), with additional terms to model uncertainty.

114 As in Eqns. (3), the process captures two sources of uncertainty in local explanations: 1) **feature**  
 115 **importance uncertainty**: the uncertainty associated with the feature importances  $\phi$ , and (2) **error**  
 116 **uncertainty**: the uncertainty associated with the error term  $\epsilon$  which captures how well our explanation  
 117  $\phi$  models the local decision surface of the underlying black box.

118 **Inference** Our inference process involves estimating the values of two key parameters:  $\phi$  and  $\sigma^2$ . By  
 119 doing so, we can compute the local explanation as well as the uncertainties associated with feature  
 120 importances and the error term. Posterior distributions on  $\phi$  and  $\sigma^2$  are normal and scaled Inv- $\chi^2$ ,  
 121 respectively, due to the corresponding conjugate priors [22]:

$$\sigma^2 | \mathcal{Z}, Y \sim \text{Scaled-Inv-}\chi^2 \left( n_0 + N, \frac{n_0 \sigma_0^2 + N s^2}{n_0 + N} \right) \quad \phi | \sigma^2, \mathcal{Z}, Y \sim \text{Normal}(\hat{\phi}, V_\phi \sigma^2) \quad (4)$$

122 Further,  $\hat{\phi}$ ,  $V_\phi$ , and  $s^2$  can be directly computed:

$$\hat{\phi} = V_\phi (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) Y) \quad V_\phi = (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + \mathbb{I})^{-1} \quad (5)$$

$$s^2 = \frac{1}{N} \left[ (Y - \mathcal{Z} \hat{\phi})^T \text{diag}(\Pi_x(\mathcal{Z})) (Y - \mathcal{Z} \hat{\phi}) + \hat{\phi}^T \hat{\phi} \right] \quad (6)$$

123 Details of the complete inference procedure including derivations of Eqns. (4-6) are provided in the  
 124 Appendix A. Note that our estimate of the posterior mean feature importances  $\hat{\phi}$  (Eqn. (5)) is the  
 125 same as that of the feature importances computed in case of LIME and KernelSHAP (Eqn. (2)).

126 **Remark 3.1.** *If we use the same proximity function  $\pi_x(z)$  in our framework as in LIME or Ker-*  
 127 *nelSHAP, the posterior mean of the feature importance  $\hat{\phi}$  output by our framework (Eq (5)) will be*  
 128 *equivalent to the feature importances output by LIME or KernelSHAP, respectively.*

129 **Feature Importance Uncertainty** To obtain the local feature importances and their associated  
 130 uncertainty, we first compute the posterior mean of the local feature importances  $\hat{\phi}$  using the closed  
 131 form expression in Eqn. (6). We then estimate the credible interval (measure of uncertainty) around  
 132 the mean feature importances by repeatedly sampling from the posterior distribution of  $\phi$  (Eq (4)).

133 **Error Uncertainty** The error term  $\epsilon$  can serve as a proxy for explanation quality because it captures  
 134 the mismatch between the constructed explanation and the local decision surface of the underlying  
 135 black box. To compute the uncertainty associated with this error term (error uncertainty), we need to  
 136 estimate the probability that the error term becomes 0 i.e.,  $P(\epsilon = 0)$ . To this end, we first calculate  
 137 the marginal posterior distribution of  $\epsilon$  by leveraging Eqn (3) and integrating out  $\sigma^2$ . This results in a  
 138 three parameter Student’s t distribution (derivation in appendix A):

$$\epsilon | \mathcal{Z}, Y \sim t_{(\nu=n_0+N)} \left( 0, \frac{n_0 \sigma_0^2 + N s^2}{n_0 + N} \right). \quad (7)$$

139 We then evaluate the probability density function (PDF) of the above posterior at 0, i.e.,  $P(\epsilon = 0)$  by  
 140 substituting the value of  $s^2$  computed using Eqn. (6) into the Student’s t distribution above (Eqn. (7)).  
 141 The resulting expression gives us the probability density that the explanation output by our framework  
 142 perfectly captures the local decision surface underlying the black box. This operation is performed in  
 143 constant time, adding minimal overhead to non-Bayesian LIME and SHAP. We illustrate how these  
 144 computed intervals capture the variance in the explanations in Figure 10.

145 **Proposition 3.2.** *As the number of perturbations around  $x$  goes to  $\infty$  i.e.,  $N \rightarrow \infty$ : (1) the estimate*  
 146 *of  $\phi$  converges to the true feature importance scores, and its uncertainty to 0. (2) uncertainty of the*  
 147 *error term  $\epsilon$  converges to the bias of the local linear model  $\phi$ . [Details in Appendix B]*

148 **BayesLIME and BayesSHAP** Our framework can be instantiated to obtain the Bayesian version of  
 149 LIME by setting the proximity function to  $\pi_x(z) = \exp(-D(x, z)^2 / \sigma^2)$  where  $D$  is a distance metric  
 150 (e.g. cosine or  $l_2$  distance), and  $n_0$  and  $\sigma_0^2$  to small values ( $10^{-6}$ ) so that the prior is uninformative.  
 151 We compute feature importance uncertainty and error uncertainty for LIME’s feature importances.

152 Our framework can also be instantiated to obtain the Bayesian version of KernelSHAP by setting  
 153 uninformative prior on  $\sigma^2$  and  $\pi_x(z) = \frac{d-1}{\binom{d}{|z|} |z| (d-|z|)}$  where  $|z|$  denotes the number of the  
 154 variables in the variable combination represented by the data point  $z$  i.e., the number of non-zero

155 valued features in the vector representation of  $z$ . Note that the original SHAP method views the  
 156 problem of constructing a local linear model as estimating the Shapley values corresponding to each  
 157 of the features [4]. These Shapley values represent the contribution of each of the features to the  
 158 black box prediction i.e.,  $f(x) = \phi_0 + \sum \phi_i$ . Therefore, the measures of uncertainty output by our  
 159 method BayesSHAP capture the reliability of the estimated variable contributions.

160 To encourage BayesLIME and BayesSHAP explanations to be sparse, we can use dimensionality  
 161 reduction or feature selection techniques as used by LIME and SHAP to obtain the top K features [2,  
 162 4, 23]. We can then construct our explanations using the data corresponding to these top K features.

### 163 3.2 Estimating the Number of Perturbations

164 One of the major drawbacks of approaches such as LIME and KernelSHAP is that they do not  
 165 provide any guidance on how to choose the number of perturbations, a key factor in obtaining reliable  
 166 explanations in an efficient manner. To address this, we leverage the uncertainty estimates output by  
 167 our framework to compute *perturbations-to-go* ( $G$ ), an estimate of how many *more* perturbations are  
 168 required to obtain explanations that satisfy a desired level of certainty. This estimate thus *predicts*  
 169 the computational cost of generating an explanation with a desired level of certainty and can help  
 170 determine whether it is even worthwhile to do so. The user specifies the confidence level of the  
 171 credible interval (denoted as  $\alpha$ ) and the *maximum* width of the credible interval ( $W$ ), e.g. “width of  
 172 95% credible interval should be less than 0.1” corresponds to  $\alpha = 0.95$  and  $W = 0.1$ . To estimate  $G$   
 173 for the local explanation of a data point  $x$ , we first generate  $S$  perturbations around  $x$  (where  $S$  is  
 174 small and chosen by the user) and fit a local linear model using our method<sup>1</sup>. This provides initial  
 175 estimates of various parameters shown in Eqns (4)-(6) which can then be used to compute  $G$ .

176 **Theorem 3.3.** *Given  $S$  seed perturbations, the number of additional perturbations required ( $G$ )  
 177 to achieve a credible interval width  $W$  of feature importance for a data point  $x$  at user-specified  
 178 confidence level  $\alpha$  can be computed as:*

$$G(W, \alpha, x) = \frac{4s_S^2}{\bar{\pi}_S \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} - S \quad (8)$$

179 where  $\bar{\pi}_S$  is the average proximity  $\pi_x(z)$  for the  $S$  perturbations,  $s_S^2$  is the empirical sum of squared  
 180 errors (SSE) between the black box and local linear model predictions, weighted by  $\pi_x(z)$ , as in (6),  
 181 and  $\Phi^{-1}(\alpha)$  is the two-tailed inverse normal CDF at confidence level  $\alpha$ .

182 *Proof (Sketch).* To estimate  $G$ , we first relate  $W$  and  $\alpha$  to  $\text{Var}(\phi_i)$ , the marginal variance of the  
 183 feature importance<sup>2</sup> for any feature  $i$ , obtained by integrating out  $\sigma^2$ . Because Student’s t can be  
 184 approximated by a Normal distribution for large degrees of freedom (here,  $S$  should be large enough),  
 185 we use the inverse normal CDF to calculate credible interval width at level  $\alpha$ . We compute  $V_\phi$  from (5)  
 186 using  $\mathcal{Z}$ , treating its entries as Bernoulli distributed with probability 0.5. Due to the covariance  
 187 structure of this sampling procedure, the resulting variance estimate after  $N$  samples is the sample  
 188 SSE  $s_S^2$  scaled by  $\approx \frac{4}{\bar{\pi}_S N}$  (derivation in appendix B). If we assume SSE scales linearly with  $S$ , we  
 189 can take this to be a reasonable estimate of  $s_N^2$  at any  $N$ . We can then estimate  $G$  as

$$\left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2 = \text{Var}(\phi_i) = \frac{4s_S^2}{\bar{\pi}_S \times (G + S)} \implies G = \frac{4s_S^2}{\bar{\pi}_S \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} - S. \quad (9)$$

190

### 191 3.3 Focused Sampling of Perturbations

192 Perturbations-to-go ( $G$ ) provides us with an estimate of how many samples are required to achieve  
 193 reliable explanations. However, if  $G$  is large, querying the black-box model for its predictions on a  
 194 large number of perturbations can be computationally expensive for larger models [24, 25]. To reduce  
 195 this cost, we develop an alternative sampling procedure called *focused sampling* which leverages

<sup>1</sup>We assume a simplified feature space where features are present or absent according to Bernoulli(.5). As in Ribeiro et al. [2], these *interpretable* features are flexible and can encode what is important to the end user.

<sup>2</sup>Since the error depends primarily on the number of perturbations,  $\text{Var}(\phi_i)$  is similar across features.

Data set	BayesLIME	BayesSHAP
ImageNet	94.8	89.9
MNIST	97.2	90.1
COMPAS	95.5	87.9
German Credit	96.9	89.6

Table 1: **Evaluating Credible Intervals.** We report the % of time the 95% credible intervals with 100 perturbations include their true values (estimated on 10,000 perturbations). Closer to 95.0 is better. Both BayesLIME and BayesSHAP are well calibrated.

uncertainty estimates to query the black box in a more targeted fashion (instead of querying randomly), thereby reducing the computational cost associated with generating reliable explanations. Inspired by active learning [26], focused sampling strategically prioritizes perturbations whose predictions the explanation is most uncertain about, when querying the black box. This enables the focused sampling procedure to query the black box only for the predictions of the most informative perturbations and thereby learn an accurate explanation with far fewer queries to the black box.

To determine how uncertain our explanation  $\phi$  is about the black box label for any given instance  $z$ , we first compute the posterior predictive distribution for  $z$  (derivation in Appendix A), given as  $\hat{y}(z)|\mathcal{Z}, Y \sim t_{(\nu=N)}(\hat{\phi}^T z, (z^T V_\phi z + 1)s^2)$ . The variance of this three parameter student’s  $t$  distribution is  $((z^T V_\phi z + 1)s^2)(N/(N - 2))$ . We refer to this variance as the *predictive variance* and it captures how uncertain our explanation  $\phi$  is about the black box label for an instance  $z$ .

The focus sampling procedure first fits the explanation with an initial  $S$  perturbations (where  $S$  is a small number). We then iterate the following procedure until the desired explanation certainty level is reached. We draw a batch of  $A$  candidate perturbations, compute their predictive variance with the Bayesian explanation, and induce a distribution over the perturbations by running softmax on the variances. We draw a batch of  $B$  perturbations from this distribution and query the black box model for their labels. Finally, we refit the Bayesian explanation on all the labeled perturbations collected so far. We provide pseudocode for the uncertainty sampling procedure in appendix Algorithm 1.

## 4 Experiments

We evaluate the proposed framework by first analyzing the quality of our uncertainty estimates i.e., feature importance uncertainty and error uncertainty. We also assess our estimates of required perturbations ( $G$ ), and evaluate the computational efficiency of focused sampling. Last, we describe a user study with 31 subjects to assess the informativeness of the explanations output by our framework.

**Setup** We experiment with a variety of real world datasets spanning multiple applications (e.g., criminal justice, credit scoring) as well as modalities (e.g., structured data, images). Our first structured dataset is **COMPAS** [27], containing criminal history, jail and prison time, and demographic attributes of 6172 defendants, with class labels that represent whether each defendant was rearrested within 2 years of release. The second structured dataset is the **German Credit** dataset from the UCI repository [28] containing financial and demographic information (including account information, credit history, employment, gender) for 1000 loan applications, each labeled as a “good” or “bad” customer. We create 80/20 train/test splits for these two datasets, and train a random forest classifier (sklearn implementation with 100 estimators) as *black box* models for each (test accuracy of 62.5% and 64.0%, respectively). We also include popular image datasets—MNIST and Imagenet. For the **MNIST** [29] handwritten digits dataset, we train a 2-layer CNN to predict the digits (test accuracy of 99.2%) and use the prediction of digit “4” as the target class. For **Imagenet** [30], we use the off-the-shelf VGG16 model [31] as the black box, and select a sample of 100 “French bulldog” images as our test set and explanation target (the model predicts French bulldog on 88% of these images). For generating explanations, we use standard implementations of the baselines LIME and KernelSHAP with default settings [2, 4]. For images, we construct super pixels as described in [2] and use them as features (number of super pixels is fixed to 20 per image). For our framework, we set perturbation size  $N = 50$ , batch size  $B = 10$ , the desired level of certainty is expressed as the width of the 95% credible interval, and use all the features.

**Quality of Uncertainty Estimates** A critical component of our explanations is the feature importance uncertainty. To evaluate the correctness of these estimates, we compute how often *true* feature importances lie within the 95% credible intervals estimated by BayesLIME and BayesSHAP. We

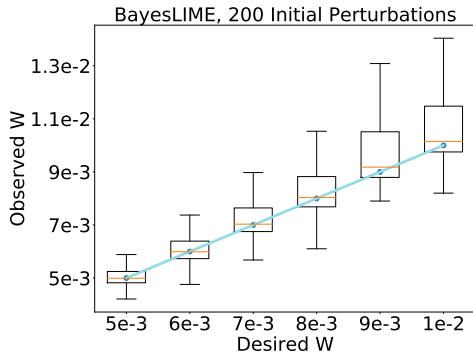


Figure 2: **Perturbations-to-go** ( $G$ ). We generate explanation with  $G$  perturbations, where  $G$  is computed using the *desired* credible interval width (x-axis), and compare desired levels to the *observed* credible interval width (y-axis) (blue line indicates ideal calibration). Results are averaged over 100 MNIST images, and  $G$ , estimated from  $S = 200$  seed samples, varies between 200 and 20,000 across images. We see that  $G$  provides a good approximation of the additional perturbations needed.

241 evaluate the quality of our credible interval estimates by running our methods with 100 perturbations  
 242 to estimate feature importances and taking the corresponding 95% credible intervals for each test  
 243 instance. We compute what fraction of the true feature importances fall within our 95% credible  
 244 intervals. Note, because there are no methods to provide uncertainty estimates for LIME and SHAP,  
 245 we do not provide further baselines. Since we do not have access to the true feature importances of  
 246 the complex black box models, following Prop 3.2, we use feature importances computed using a  
 247 large value of  $N$  ( $N = 10,000$ ), and treat the resulting estimates as ground truth.

248 Results for BayesLIME in Table 1 indicate that the true feature importances are close to ideal and  
 249 indicate the estimates are well calibrated. While the estimates by BayesSHAP are somewhat less  
 250 calibrated (true feature importances fall within our estimated 95% credible intervals about 89.6 to  
 251 90.1% of the time), they still are quite close to ideal. All in all, these results confirm that the credible  
 252 intervals learned by our methods are well calibrated and therefore highly reliable in capturing the  
 253 uncertainty of the feature importances. Lastly, though we set our priors to be uninformative in general,  
 254 we also investigate how sensitive our uncertainty estimates are to hyperparameter choices in Figure 5  
 255 in the Appendix. We find that the explanation uncertainty becomes uncalibrated with strong priors.  
 256 However, our explanations seem to be robust to hyperparameter choices in general.

257 **Correctness of Estimated Number of Perturbations** We assess whether our estimate of  
 258 *perturbations-to-go* ( $G$ ; Section 3.2) is an accurate estimate of the *additional* number of pertur-  
 259 bations needed to reach a desired level of feature importance certainty. We carry out this experiment  
 260 on MNIST data (additional datasets explored in Appendix C) and use  $S = 200$  as the initial number  
 261 of perturbations to obtain a preliminary explanation and its associated uncertainty estimates. We then  
 262 leverage these estimates to compute  $G$  for 6 different certainty levels. First, we observe significant  
 263 differences in  $G$  estimates across instances (details in appendix C) i.e. number of perturbations  
 264 needed to obtain a particular level of certainty varied significantly across instances—ranging from  
 265 200-5,000 for the lowest level of certainty to 200-20,000 for higher levels of certainty. Next, for  
 266 each image and certainty level, we run our method for the estimated number of perturbations ( $G$ )  
 267 to determine if the observed estimates of uncertainty (observed credible interval width  $W$ ) match  
 268 the desired levels of uncertainty (desired credible interval width  $W$ ). Results in Figure 2 show that  
 269 the observed and desired levels of certainty are well calibrated, demonstrating that  $G$  estimates are  
 270 reliable approximations of the additional number of perturbations needed.

271 **Efficiency of Focused Sampling** *Focused sampling* uses the *predictive variance* to strategically  
 272 choose perturbations that will reduce uncertainty in order to be labeled by the black box (section 3.3).  
 273 Here, we will evaluate the efficiency of the focused sampling procedure. First, we assess whether  
 274 focused sampling produces reliable explanations (as measured by error uncertainty ( $P(\epsilon = 0)$ )) more  
 275 efficiently than random sampling. To this end, we experiment with BayesLIME on Imagenet data to  
 276 carry out this analysis. This setting replicates scenarios where LIME is applied to a computationally  
 277 expensive black box model, making it highly desirable to limit the number of perturbations to reduce  
 278 total running time. We run each sampling strategy for 2,000 perturbations and plot the number of  
 279 model queries versus error uncertainty. The results in Figure 3 show that focused sampling results in  
 280 faster convergence to reliable and high quality explanations; focused sampling stabilizes within a  
 281 couple hundred model queries while random sampling takes over 1,000. Note, as the inefficiency of  
 282 querying the black box model increases, the advantages of focused sampling decreasing total running

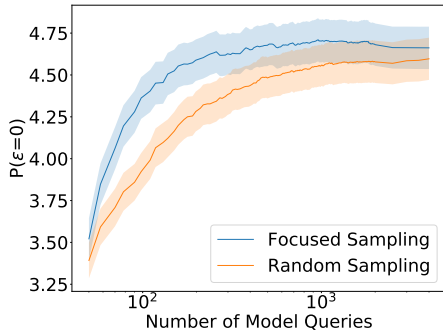


Figure 3: **Efficiency of focused sampling** for 100 Imagenet images, with random sampling as a baseline. We provide mean and standard error. We assess the efficiency of focused sampling by comparing *error uncertainty* (y-axis) over model queries, and show much quicker convergence than random sampling.

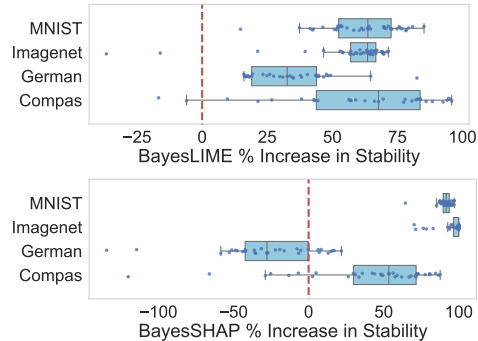


Figure 4: **Assessing the % increase in stability** of BayesLIME and BayesSHAP over LIME and SHAP respectively. Our Bayesian methods are significant more stable ( $\rho < 1e-2$  according to Wilcoxon signed-rank test) except for BayesSHAP on German Credit, where there is not a significant difference between the methods ( $\rho > 0.05$ ).

283 time of the explanations will only become more pronounced. These results clearly demonstrate that  
 284 focused sampling can significantly speed up the process of generating high quality local explanations.

285 We also benchmark the efficiency of BayesLIME and BayesSHAP against Guo et al. [32], a related  
 286 Bayesian explanation method that uses a Bayesian non parametric mixture regression and MCMC  
 287 for parameter inference. Fixing their mixture regression to a single component results in a similar  
 288 model to ours and thus is a useful point of comparison. To explain a single instance on ImageNet  
 289 using VGG16, their approach takes 139.2 seconds, while BayesLIME and BayesSHAP take 20.3  
 290 seconds and 21.1 seconds respectively, under the same conditions, demonstrating that the closed form  
 291 solution is very efficient. Additionally, in Appendix C, we also check if focused sampling causes any  
 292 bias (due to sampling based on uncertainty estimates) that results in convergence to a different/wrong  
 293 explanation, however our results clearly indicate that this is not the case.

294 **Stability of BayesLIME & BayesSHAP** Recall that LIME & SHAP are not stable: small changes  
 295 to instances can produce substantially different explanations. We consider whether BayesLIME &  
 296 BayesSHAP produce more stable explanations than their LIME & SHAP counterparts. To perform  
 297 this analysis, we use the local Lipschitz metric for explanation stability [18]:

$$\hat{L}(x_i) = \operatorname{argmax}_{x_j \in N_\epsilon(x_i)} \frac{\|\phi_i - \phi_j\|_2}{\|x_i - x_j\|_2} \quad (10)$$

298 where  $x_i$  refers to an instance,  $N_\epsilon(x_i)$  is the  $\epsilon$ -ball centered at  $x_i$ , and  $\phi_i$  and  $\phi_j$  are the explanation  
 299 parameters for  $x_i$  and  $x_j$ . Lower values of the above metric indicate more stable explanations. We  
 300 follow the setup outline by Alvarez-Melis and Jaakkola [18] and compute the local Lipschitz values,  
 301 comparing both LIME & BayesLIME and SHAP & BayesSHAP across Compas, German Credit,  
 302 MNIST, and Imagenet. We perform the comparison using the default number of perturbations in  
 303 both LIME & SHAP, and use this same number in the respective Bayesian variants. We use focused  
 304 sampling for BayesLIME and BayesSHAP, and report the % increase in stability of these approaches  
 305 over LIME and SHAP, respectively, for 40 points in the test data. The results given in Figure 4 show  
 306 a clear improvement (on average 53%) in stability for both BayesLIME and BayesSHAP in all cases  
 307 except German Credit for BayesSHAP. Further, we run a Wilcoxon signed-rank test and find our  
 308 results are statistically significant in all cases ( $\rho < 1e-2$ ) except for BayesSHAP for German Credit,  
 309 where there is not a significant difference between the methods ( $\rho > 0.05$ ). These results demonstrate  
 310 BayesLIME and BayesSHAP are more stable than previous methods.

311 **User Study** We perform a user study with 31 subjects to compare BayesLIME and LIME explanations  
 312 on MNIST. We evaluate the following: are explanations with low levels of uncertainty (i.e., most  
 313 confident explanations) more meaningful to humans? To answer this question, we mask the most  
 314 important features selected by BayesLIME and LIME, and ask users to guess the digit of the masked  
 315 images. The better the explanation, the more difficult it should be for the users to get it right. We

316 randomly select 15 correctly predicted test images, generate explanations by sweeping over a range  
317 of perturbation amounts  $[10^{-5}, \dots, 10^{3.5}]$  incremented by 0.5. We choose the *top* explanation for each  
318 image based on either fidelity (for LIME) or  $P(\epsilon = 0)$  (for BayesLIME). We sent the user study out  
319 to students and researchers with background in computer science. A screen shot of the task is show  
320 in Figure 8 in the Appendix. We find that the explanations output by our methods focus on more  
321 informative parts of the image, since hiding them makes it difficult for humans to guess the digit.  
322 Users had an error rate of 25.7% for LIME, while it was 30.7% for BayesLIME, both with standard  
323 error 0.003 ( $\rho = 0.028$  through a one-tailed two sample t-test). This result indicates that our method  
324 BayesLIME and the associated measure of explanation uncertainty result in more high quality and  
325 reliable explanations compared to LIME and its associated fidelity metric.

## 326 5 Related Work

327 **Interpretability Methods** A variety of interpretability methods have been proposed. Some methods  
328 that are inherently interpretable include additive models [33, 34], decision lists and sets [35, 36], and  
329 instance-based explanations [37]. However, black-box models are often more flexible, accurate, and  
330 easier to use; thus, there has been a lot of interest in constructing post hoc explanations[38]. These  
331 include LIME [2] and SHAP [4], which are among the most popular due to their broad applicability  
332 and code availability, but saliency maps [5–8], permutation feature importance [39], and partial  
333 dependency plots [40] also follow this paradigm. Other approaches to post hoc explanations focus on  
334 rule-based models [1, 3], counterfactuals [41, 42], and influence functions [9].

335 **Vulnerabilities of Post hoc Explanations** Recent work has shed light on the downsides of post hoc  
336 explanation techniques. These methods are often highly sensitive to small changes in inputs [14],  
337 are susceptible to manipulation [15, 16, 43, 44], and are not faithful to the underlying black boxes  
338 [45]. Perturbation-based explanation methods such as LIME and SHAP are subject to additional  
339 criticisms: results vary between runs of the algorithms [18–20, 46, 21], and hyperparameters used to  
340 select the perturbations can greatly influence the resulting explanation [20]. Prior work has attempted  
341 to tackle the problem of instability in perturbation-based explanations by averaging over several  
342 explanations [47, 19], however, this is computationally expensive. Other works related to creating  
343 more trustworthy explanations include development of sanity checks for explainers [48, 17, 49].  
344 These techniques represent an important step towards improved usability, given experimental evidence  
345 that humans are often too eager to accept inaccurate machine explanations [50–53]. Recent works  
346 theoretically analyze the sources of non-robustness in black box explanations [54–56].

347 **Bayesian Methods in Explainable ML** Few recent works have adopted Bayesian formulations to  
348 explain black box models [32, 57]. Guo et al. [32] introduce a Bayesian non-parametric approach to  
349 fit a *global* surrogate model. Their formulation seeks to fit a mixture of generalizable explanations  
350 across instances. Zhao et al. [57] draw on Bayesian frameworks from a preprint of this work to  
351 study whether incorporating informative priors improves the stability of the resulting explanations.  
352 However, neither of these works focus on modeling the uncertainty of local explanations. Further,  
353 these approaches also do not tackle the critical problems of estimating key hyperparameters or  
354 improving efficiency of computing explanations.

## 355 6 Conclusion

356 We developed a Bayesian framework for generating local explanations along with their associated  
357 uncertainty. We instantiated this framework to obtain Bayesian versions of LIME and SHAP that  
358 output pointwise estimates of feature importances as well as their associated credible intervals. These  
359 intervals enabled us to infer the quality of the explanations and output explanations that satisfied user  
360 specified levels of uncertainty. We carried out theoretical analysis that leverages these uncertainty  
361 measures (credible intervals) to estimate the values of critical hyperparameters (e.g., the number of  
362 perturbations). We also proposed a novel sampling technique called focused sampling that leverages  
363 uncertainty estimates to determine how to sample perturbations for faster convergence. One potential  
364 negative outcome of our work is that if the credible intervals output by our framework are incorrect, it  
365 could lead to overreliance of faulty explanations. It would be interesting to extend our framework to  
366 produce global explanations with uncertainty guarantees and explore how uncertainty quantification  
367 can help calibrate user trust in model explanations.

368 **References**

- 369 [1] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-  
370 agnostic explanations. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- 371 [2] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why Should I Trust You? explaining  
372 the predictions of any classifier. In *Knowledge Discovery and Data mining (KDD)*, 2016.
- 373 [3] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. Faithful and customizable  
374 explanations of black box models. In *Proceedings of the 2019 AAAI/ACM Conference on AI,  
375 Ethics, and Society*, pages 131–138. ACM, 2019.
- 376 [4] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In  
377 *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- 378 [5] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:  
379 Visualising image classification models and saliency maps. In *Workshop at International  
380 Conference on Learning Representations*, 2014.
- 381 [6] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks.  
382 In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages  
383 3319–3328. JMLR. org, 2017.
- 384 [7] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi  
385 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based  
386 localization. In *ICCV*, 2017.
- 387 [8] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smooth-  
388 grad: removing noise by adding noise. *Workshop on Visualization for Deep Learning, ICML,*  
389 *2017*.
- 390 [9] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions.  
391 In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages  
392 1885–1894. JMLR. org, 2017.
- 393 [10] Osbert Bastani, Carolyn Kim, and Hamsa Bastani. Interpretability via model extraction. *FAT/ML  
394 Workshop 2017*, 2017.
- 395 [11] Gregory Plumb, Denali Molitor, and Ameet S Talwalkar. Model agnostic supervised local  
396 explanations. In *Neural Information Processing Systems*, 2018.
- 397 [12] Radwa Elshawi, Mouaz H Al-Mallah, and Sherif Sakr. On the interpretability of machine  
398 learning-based model for predicting hypertension. *BMC medical informatics and decision  
399 making*, 19(1):146, 2019.
- 400 [13] Leanne S Whitmore, Anthe George, and Corey M Hudson. Mapping chemical performance on  
401 molecular structures using locally interpretable explanations. *arXiv preprint arXiv:1611.07443*,  
402 2016.
- 403 [14] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile.  
404 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3681–3688,  
405 2019.
- 406 [15] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. Fooling  
407 lime and shap: Adversarial attacks on post hoc explanation methods. *Conference on Artificial  
408 Intelligence, Ethics, and Society (AIES)*, 2020.
- 409 [16] Ann-Kathrin Dombrowski, Maximilian Alber, Christopher J Anders, Marcel Ackermann, Klaus-  
410 Robert Müller, and Pan Kessel. Explanations can be manipulated and geometry is to blame.  
411 *arXiv preprint arXiv:1906.07983*, 2019.
- 412 [17] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim.  
413 Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems*, pages  
414 9505–9515, 2018.

- 415 [18] David Alvarez-Melis and Tommi S. Jaakkola. On the robustness of interpretability methods.  
416 *ICML Workshop on Human Interpretability in Machine Learning*, 2018.
- 417 [19] Eunjin Lee, David Braines, Mitchell Stiffler, Adam Hudler, and Daniel Harborne. Developing  
418 the sensitivity of lime for better machine learning explanation. In *Artificial Intelligence and*  
419 *Machine Learning for Multi-Domain Operations Applications*, volume 11006, page 1100610.  
420 International Society for Optics and Photonics, 2019.
- 421 [20] Hui Fen Tan, Kuangyan Song, Madeilene Udell, Yiming Sun, and Yujia Zhang. “why should  
422 you trust my explanation?” understanding uncertainty in lime explanations. In *ICML Workshop*  
423 *on AI for Social Good*, 2019.
- 424 [21] Jianbo Chen, Le Song, Martin J. Wainwright, and Michael I. Jordan. L-shapley and c-shapley:  
425 Efficient model interpretation for structured data. In *International Conference on Learning*  
426 *Representations*, 2019.
- 427 [22] Andrew Moore. Locally weighted bayesian regression, January 1995.
- 428 [23] Kacper Sokol, Alexander Hepburn, Raúl Santos-Rodríguez, and Peter A. Flach. blimey:  
429 Surrogate prediction explanations beyond lime. *NeurIPS HCML Workshop*, 2019.
- 430 [24] Emily L. Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting  
431 linear structure within convolutional networks for efficient evaluation. In *NIPS*, 2014.
- 432 [25] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural  
433 networks with low rank expansions. *BMVC 2014 - Proceedings of the British Machine Vision*  
434 *Conference 2014*, 05 2014.
- 435 [26] Burr Settles. Active learning literature survey. 2010.
- 436 [27] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias. In *ProPublica*,  
437 2016.
- 438 [28] Dheeru Dua and Casey Graff. Uci machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- 440 [29] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs*  
441 *[Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 442 [30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale  
443 Hierarchical Image Database. In *CVPR09*, 2009.
- 444 [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
445 image recognition. In *ICLR*, 2015.
- 446 [32] Wenbo Guo, Sui Huang, Yunzhe Tao, Xinyu Xing, and Lin Lin. Explaining deep learning  
447 models – a bayesian non-parametric approach. In *Neural Information Processing Systems*  
448 *(NeurIPS)*. 2018.
- 449 [33] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. Accurate intelligible models with  
450 pairwise interactions. In *KDD*, 2013.
- 451 [34] Berk Ustun, Stefano Traca, and Cynthia Rudin. Supersparse linear integer models for inter-  
452 pretable classification. *arXiv preprint arXiv:1306.6677*, 2013.
- 453 [35] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. Interpretable decision sets: A joint  
454 framework for description and prediction. In *Knowledge Discovery and Data mining (KDD)*,  
455 2016.
- 456 [36] Elaine Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin.  
457 Learning certifiably optimal rule lists for categorical data. *arXiv preprint arXiv:1704.01701*,  
458 2017.
- 459 [37] Been Kim, Cynthia Rudin, and Julie Shah. The bayesian case model: A generative approach for  
460 case-based reasoning and prototype classification. *arXiv preprint arXiv:1503.01161*, 2015.

- 461 [38] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-Agnostic Interpretability of  
462 Machine Learning. In *ICML Workshop on Human Interpretability in Machine Learning (WHI)*,  
463 June 2016.
- 464 [39] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 465 [40] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of*  
466 *statistics*, pages 1189–1232, 2001.
- 467 [41] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In  
468 *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT\* ’19*, pages  
469 10–19, 2019. ISBN 978-1-4503-6125-5.
- 470 [42] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without  
471 opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.
- 472 [43] Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via  
473 adversarial model manipulation. In *Advances in Neural Information Processing Systems 32*,  
474 pages 2921–2932. 2019.
- 475 [44] Junlin Wang, Jens Tuyls, Eric Wallace, and Sameer Singh. Gradient-based Analysis of NLP  
476 Models is Manipulable. In *Findings of the Association for Computational Linguistics: EMNLP*  
477 *(EMNLP Findings)*, page 247–258, 2020.
- 478 [45] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions  
479 and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206, 2019.
- 480 [46] Muhammad Rehman Zafar and Naimul Mefraz Khan. Dlime: a deterministic local interpretable  
481 model-agnostic explanations approach for computer-aided diagnosis systems. *arXiv preprint*  
482 *arXiv:1906.10263*, 2019.
- 483 [47] Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On  
484 the (in) fidelity and sensitivity of explanations. In *Advances in Neural Information Processing*  
485 *Systems*, pages 10965–10976, 2019.
- 486 [48] Oana-Maria Camburu, Eleonora Giunchiglia, Jakob Foerster, Thomas Lukasiewicz, and Phil  
487 Blunsom. Can i trust the explainer? verifying post-hoc explanatory methods. *arXiv preprint*  
488 *arXiv:1910.02065*, 2019.
- 489 [49] Mengjiao Yang and Been Kim. Bim: Towards quantitative evaluation of interpretability methods  
490 with ground truth. *arXiv:1907.09701*, 2019.
- 491 [50] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jen-  
492 nifer Wortman Vaughan. Interpreting interpretability: Understanding data scientists’ use of  
493 interpretability tools for machine learning. In *CHI*, April 2020.
- 494 [51] Fred Hohman, Andrew Head, Rich Caruana, Robert DeLine, and Steven M Drucker. Gamut:  
495 A design probe to understand how data scientists understand machine learning models. In  
496 *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13,  
497 2019.
- 498 [52] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan,  
499 and Hanna Wallach. Manipulating and measuring model interpretability. *arXiv preprint*  
500 *arXiv:1802.07810*, 2018.
- 501 [53] Himabindu Lakkaraju and Osbert Bastani. " how do i fool you?" manipulating user trust via  
502 misleading black box explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics,*  
503 *and Society*, pages 79–85, 2020.
- 504 [54] Damien Garreau and Ulrike von Luxburg. Looking deeper into lime. *arXiv preprint*  
505 *arXiv:2008.11092*, 2020.
- 506 [55] Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise  
507 explanations of neural networks using adversarial training. In *International Conference on*  
508 *Machine Learning*, pages 1383–1391. PMLR, 2020.

- 509 [56] Alexander Levine, Sahil Singla, and Soheil Feizi. Certifiably robust interpretation in deep  
510 learning. *arXiv preprint arXiv:1905.12105*, 2019.
- 511 [57] Xingyu Zhao, Xiaowei Huang, Valentin Robu, and David Flynn. Baylime: Bayesian local  
512 interpretable model-agnostic explanations. *arXiv preprint arXiv:2012.03058*, 2020.
- 513 [58] Ludwig Fahrmeir, Thomas Kneib, and Stefan Lang. *Regression*. Statistik und ihre Anwendungen.  
514 Springer, 2007. ISBN 978-3-540-33932-8.
- 515 [59] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.

## 516 Checklist

- 517 1. For all authors...
- 518 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
519 contributions and scope? [Yes] We introduce our Bayesian explanation framework and  
520 describe that it leads to improved explanations. By describing in the main text each  
521 way it does this (e.g., focused sampling, PTG) and validating these experimentally, we  
522 validate the claims made in the introduction.
- 523 (b) Did you describe the limitations of your work? [Yes] We describe the limitations of our  
524 work in a number of places. For one, we describe where our Bayesian explanations are  
525 not successful experimentally (e.g., German credit in BayesSHAP stability in figure 4).  
526 We also indicate the need to study how humans perceive the uncertainty estimates in  
527 the Conclusion.
- 528 (c) Did you discuss any potential negative societal impacts of your work? [Yes] We discuss  
529 potential impacts in the Conclusion.
- 530 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
531 them? [Yes]
- 532 2. If you are including theoretical results...
- 533 (a) Did you state the full set of assumptions of all theoretical results? [Yes] The assump-  
534 tions for our main theorem are clearly stated in the theorem body.
- 535 (b) Did you include complete proofs of all theoretical results? [Yes] We provide proofs in  
536 the appendix.
- 537 3. If you ran experiments...
- 538 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
539 imental results (either in the supplemental material or as a URL)? [Yes] Provided in  
540 supplement.
- 541 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
542 were chosen)? [Yes] We provide details in the experiment section, Section 4.
- 543 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
544 ments multiple times)? [Yes]
- 545 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
546 of GPUs, internal cluster, or cloud provider)? [Yes] We provide discussion in section D  
547 in the Appendix.
- 548 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 549 (a) If your work uses existing assets, did you cite the creators? [Yes] We use public datasets  
550 and cite the creators.
- 551 (b) Did you mention the license of the assets? [Yes] We mention licenses in Section E
- 552 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
553 We include code in the supplement.
- 554 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
555 using/curating? [Yes] The datasets used are commonly used public datasets. For our  
556 user study, we provided a page where we asked participants whether they were willing  
557 to participate.

- 558 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
559 information or offensive content? [Yes] The datasets used are commonly used public  
560 datasets. For our user study, we did not ask for personal information.
- 561 5. If you used crowdsourcing or conducted research with human subjects...
- 562 (a) Did you include the full text of instructions given to participants and screenshots, if  
563 applicable? [Yes] Provided in Appendix.
- 564 (b) Did you describe any potential participant risks, with links to Institutional Review  
565 Board (IRB) approvals, if applicable? [N/A]
- 566 (c) Did you include the estimated hourly wage paid to participants and the total amount  
567 spent on participant compensation? [N/A]