

# Reliable Post hoc Explanations: Modeling Uncertainty in Explainability

<b>Dylan Slack</b> UC Irvine dslack@uci.edu	<b>Sophie Hilgard</b> Harvard University ash798@g.harvard.edu	<b>Sameer Singh</b> UC Irvine sameer@uci.edu	<b>Himabindu Lakkaraju</b> Harvard University hlakkaraju@hbs.edu
---	---	--	--

## Abstract

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

## 1 Introduction

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

Existing local explanation methods, however, suffer from several drawbacks. Explanations generated using these methods may be unstable [14–18], i.e., negligibly small perturbations to an instance can result in substantially different explanations. These methods are also inconsistent [19] i.e., multiple runs on the same input instance with the same parameter settings may result in vastly different explanations. There are also no reliable metrics to ascertain the quality of the explanations

<sup>1</sup>Project Page: <https://dylanslacks.website/reliable/index.html>

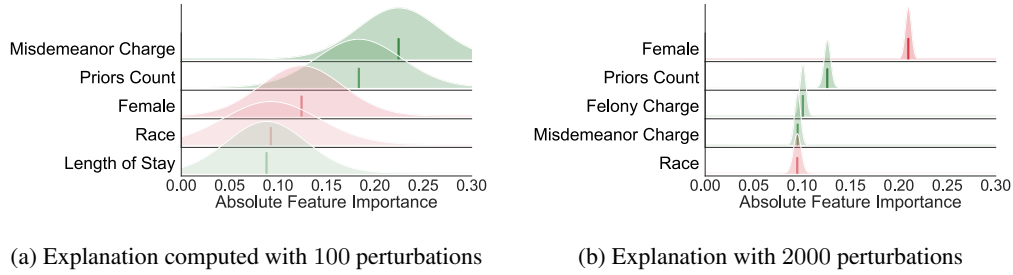


Figure 1: **Example explanations** on for an instance from the COMPAS dataset, where vertical lines indicate the feature importance by LIME (red is negative effect, green is positive) and the shaded region visualizes the uncertainty estimated by BayesLIME. While LIME produces very different and contradictory 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.

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

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

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

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

## 2 Notation & Background

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

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

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

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

The above objective function has the following closed form solution:

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

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

## 3 Our Framework: Bayesian Local Explanations

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

### 3.1 Constructing Bayesian Local Explanations

Our goal here is to explain the behavior of a given black box model  $f$  in the vicinity of an instance  $x$  while also capturing the uncertainty associated with the explanation. To this end, we propose a Bayesian framework for constructing local linear model based explanations and capturing their associated uncertainty. We model the black box prediction of each perturbation  $z$  as a linear combination of the corresponding feature values ( $\phi^T z$ ) plus an error term ( $\epsilon$ ) as shown in Eqn (4). While the weights of the linear combination  $\phi$  capture the feature importances and thereby constitute our explanation,  $\epsilon$  captures the error that arises due to the mismatch between our explanation  $\phi$  and the local decision surface of the black box model  $f$ . Our complete generative process is shown below:

$$y|z, \phi, \epsilon \sim \phi^T z + \epsilon \quad \epsilon \sim \mathcal{N}(0, \frac{\sigma^2}{\pi_x(z)}) \quad (3)$$

$$\phi|\sigma^2 \sim \mathcal{N}(0, \sigma^2 \mathbb{I}) \quad \sigma^2 \sim \text{Inv-}\chi^2(n_0, \sigma_0^2). \quad (4)$$

The error term is modeled as a Gaussian whose variance relies on the proximity function  $\pi_x(z)$  i.e.,  $\epsilon \sim \mathcal{N}(0, \frac{\sigma^2}{\pi_x(z)})$ . This proximity function ensures that perturbations closer to the data point  $x$  are modeled accurately, while allowing more room for error in case of perturbations that are farther away.  $\pi_x(z)$  can be computed using cosine or  $l_2$  distance or other game theoretic principles similar to that of LIME and KernelSHAP (see Section 2). The conjugate priors on  $\phi$  and  $\sigma^2$  are shown in Eqn (4). Note that, the distributions on error  $\epsilon$  and feature importance  $\phi$  both consider the parameter  $\sigma^2$ . The

fact that the prior on the feature importances considers  $\sigma^2$  has an intuitive interpretation: if we have prior knowledge that the error of the explanation is small, we expect to be more confident about the feature importances. Similarly, if we have prior knowledge the error is large, we expect to be less confident about the feature importances.

Thus, our generative process corresponds to the Bayesian version of the weighted least squares formulation of LIME and KernelSHAP outlined in Eqn. (1), with additional terms to model uncertainty. As in Eqns. (4), the process captures two sources of uncertainty in local explanations: 1) **feature importance uncertainty**: the uncertainty associated with the feature importances  $\phi$ , and (2) **error uncertainty**: the uncertainty associated with the error term  $\epsilon$  which captures how well our explanation  $\phi$  models the local decision surface of the underlying black box.

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

$$\begin{aligned}\sigma^2 | \mathcal{Z}, Y &\sim \text{Scaled-Inv-}\chi^2 \left( n_0 + N, \frac{n_0 \sigma_0^2 + N s^2}{n_0 + N} \right) \\ \phi | \sigma^2, \mathcal{Z}, Y &\sim \text{Normal}(\hat{\phi}, V_\phi \sigma^2)\end{aligned}\tag{5}$$

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

$$\begin{aligned}\hat{\phi} &= V_\phi (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) Y) \\ V_\phi &= (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + \mathbb{I})^{-1}\end{aligned}\tag{6}$$

$$s^2 = \frac{1}{N} \left[ (Y - \mathcal{Z} \hat{\phi})^T \text{diag}(\Pi_x(\mathcal{Z})) (Y - \mathcal{Z} \hat{\phi}) + \hat{\phi}^T \hat{\phi} \right]\tag{7}$$

Details of the complete inference procedure including derivations of Eqns. (5-7) are provided in the Appendix A. Note that our estimate of the posterior mean feature importances  $\hat{\phi}$  (Eqn. (6)) is the same as that of the feature importances computed in case of LIME and KernelSHAP (Eqn. (2)).

**Remark 3.1.** *If we use the same proximity function  $\pi_x(z)$  in our framework as in LIME or KernelSHAP, the posterior mean of the feature importance  $\hat{\phi}$  output by our framework (Eq (6)) will be equivalent to the feature importances output by LIME or KernelSHAP, respectively.*

**Feature Importance Uncertainty** To obtain the local feature importances and their associated uncertainty, we first compute the posterior mean of the local feature importances  $\hat{\phi}$  using the closed form expression in Eqn. (7). We then estimate the credible interval (measure of uncertainty) around the mean feature importances by repeatedly sampling from the posterior distribution of  $\phi$  (Eq (5)).

**Error Uncertainty** The error term  $\epsilon$  can serve as a proxy for explanation quality because it captures the mismatch between the constructed explanation and the local decision surface of the underlying black box. We first calculate the marginal posterior distribution of  $\epsilon$  by leveraging Eqn (4) and integrating out  $\sigma^2$ . This results in a three parameter Student's t distribution (derivation in appendix A):

$$\epsilon | \mathcal{Z}, Y \sim t_{(v=n_0+N)} \left( 0, \frac{n_0 \sigma_0^2 + N s^2}{n_0 + N} \right).\tag{8}$$

We then evaluate the probability density function (PDF) of the above posterior at 0, i.e.,  $P(\epsilon = 0)$  by substituting the value of  $s^2$  computed using Eqn. (7) into the Student's t distribution above (Eqn. (8)). The resulting expression gives us the probability density that the explanation output by our framework perfectly captures the local decision surface underlying the black box. This operation is performed in constant time, adding minimal overhead to non-Bayesian LIME and SHAP. We illustrate how these computed intervals capture the variance in the explanations in Figure 9.

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

**BayesLIME and BayesSHAP** Our framework can be instantiated to obtain the Bayesian version of LIME by setting the proximity function to  $\pi_x(z) = \exp(-D(x, z)^2/\sigma^2)$  where  $D$  is a distance metric

(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. We compute feature importance uncertainty and error uncertainty for LIME’s feature importances.

Our framework can also be instantiated to obtain the Bayesian version of KernelSHAP by setting uninformative prior on  $\sigma^2$  and  $\pi_x(z) = \frac{d-1}{(d \text{ choose } |z|)|z|(d-|z|)}$  where  $|z|$  denotes the number of the variables in the variable combination represented by the data point  $z$  i.e., the number of non-zero valued features in the vector representation of  $z$ . Note that the original SHAP method views the problem of constructing a local linear model as estimating the Shapley values corresponding to each of the features [4]. These Shapley values represent the contribution of each of the features to the black box prediction i.e.,  $f(x) = \phi_0 + \sum \phi_i$ . Therefore, the measures of uncertainty output by our method BayesSHAP capture the reliability of the estimated variable contributions.

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

### 3.2 Estimating the Number of Perturbations

One of the major drawbacks of approaches such as LIME and KernelSHAP is that they do not provide any guidance on how to choose the number of perturbations, a key factor in obtaining reliable explanations in an efficient manner. To address this, we leverage the uncertainty estimates output by our framework to compute *perturbations-to-go* ( $G$ ), an estimate of how many *more* perturbations are required to obtain explanations that satisfy a desired level of certainty. This estimate thus *predicts* the computational cost of generating an explanation with a desired level of certainty and can help determine whether it is even worthwhile to do so. The user specifies the confidence level of the credible interval (denoted as  $\alpha$ ) and the *maximum* width of the credible interval ( $W$ ), e.g. “width of 95% credible interval should be less than 0.1” corresponds to  $\alpha = 0.95$  and  $W = 0.1$ . To estimate  $G$  for the local explanation of a data point  $x$ , we first generate  $S$  perturbations around  $x$  (where  $S$  is small and chosen by the user) and fit a local linear model using our method<sup>2</sup>. This provides initial estimates of various parameters shown in Eqns (5)-(7) which can then be used to compute  $G$ .

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

$$G(W, \alpha, x) = \frac{4s_S^2}{\bar{\pi}_S \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} - S \quad (9)$$

where  $\bar{\pi}_S$  is the average proximity  $\pi_x(z)$  for the  $S$  perturbations,  $s_S^2$  is the empirical sum of squared errors (SSE) between the black box and local linear model predictions, weighted by  $\pi_x(z)$ , as in (7), and  $\Phi^{-1}(\alpha)$  is the two-tailed inverse normal CDF at confidence level  $\alpha$ .

*Proof (Sketch).* To estimate  $G$ , we first relate  $W$  and  $\alpha$  to  $\text{Var}(\phi_i)$ , the marginal variance of the feature importance<sup>3</sup> for any feature  $i$ , obtained by integrating out  $\sigma^2$ . Because Student’s t can be approximated by a Normal distribution for large degrees of freedom (here,  $S$  should be large enough), we use the inverse normal CDF to calculate credible interval width at level  $\alpha$ . We compute  $V_\phi$  from (6) using  $\mathcal{Z}$ , treating its entries as Bernoulli distributed with probability 0.5. Due to the covariance structure of this sampling procedure, the resulting variance estimate after  $N$  samples is the sample 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 can take this to be a reasonable estimate of  $s_N^2$  at any  $N$ . We can then estimate  $G$  as

$$\left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2 = \text{Var}(\phi_i) = \frac{4s_S^2}{\bar{\pi}_S \times (G + S)} \implies G = \frac{4s_S^2}{\bar{\pi}_S \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} - S. \quad (10)$$

□

<sup>2</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>3</sup>Since the error depends primarily on the number of perturbations,  $\text{Var}(\phi_i)$  is similar across features.

### 3.3 Focused Sampling of Perturbations

Perturbations-to-go ( $G$ ) provides us with an estimate of how many samples are required to achieve reliable explanations. However, if  $G$  is large, querying the black-box model for its predictions on a large number of perturbations can be computationally expensive for larger models [24, 25]. To reduce this cost, we develop an alternative sampling procedure called *focused sampling* which leverages 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  $\hat{y}(z)|\mathcal{Z}, Y \sim t_{(V=N)}(\hat{\phi}^T z, (z^T V_\phi z + 1)s^2)$ . The variance of this three parameter student's  $t$  distribution is,

$$\text{var}(\hat{y}(z)) = ((z^T V_\phi z + 1)s^2)(N/(N - 2)) \quad (11)$$

We refer to this variance as the *predictive variance*  $\text{var}(\hat{y}(z))$ , and it captures how uncertain our explanation  $\phi$  is about the black box prediction.

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 with temperature parameter  $\tau$ . 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 Algorithm 1.

---

#### Algorithm 1 Focused sampling for local explanations

---

**Require:** Model  $f$ , Data instance  $x$ , Number of perturbations  $N$ , Number of seed perturbations  $S$ , Batch size  $B$ , Pool size  $A$ , temperature  $\tau$

- 1: **function** FOCUSED SAMPLE
- 2:   Initialize  $\mathcal{Z}$  with  $S$  seed perturbations.
- 3:   Fit  $\hat{\phi}$  on  $\mathcal{Z}$  ▷ Using Eqn (6)
- 4:   **for**  $i \leftarrow 1$  to  $N - S$  in increments of  $B$  **do**
- 5:      $\mathcal{Q} \leftarrow$  Generate  $A$  candidate perturbations
- 6:     Compute  $\text{var}(\hat{y}(z))$  on  $\mathcal{Q}$  ▷ Using Eqn (11)
- 7:     Define  $\mathcal{Q}_{\text{dist}}$  as  $\propto \exp(\text{var}(\hat{y}(z))/\tau)$
- 8:      $\mathcal{Q}_{\text{new}} \leftarrow$  Draw  $B$  samples from  $\mathcal{Q}_{\text{dist}}$
- 9:      $\mathcal{Z} \leftarrow \mathcal{Z} \cup \mathcal{Q}_{\text{new}}$ ; Fit  $\hat{\phi}$  on  $\mathcal{Z}$  ▷ Using Eqn (6)
- 10:   **end for**
- 11:   **return**  $\hat{\phi}$
- 12: **end function**

---

## 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



	BayesLIME	BayesSHAP		BayesLIME	BayesSHAP
TABULAR DATASETS			MNIST		
COMPAS	95.5	87.9	Digit 1	95.8	98.4
German Credit	96.9	89.6	Digit 2	95.8	97.4
IMAGENET			Digit 3	95.2	96.3
Corn	94.6	91.8	Digit 4	97.2	90.1
Broccoli	91.4	89.2	Digit 5	95.2	95.6
French Bulldog	94.8	89.9	Digit 6	96.7	96.8
Scuba Diver	92.4	94.6	Digit 7	95.7	95.3

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.

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 82.8% and 72.5%, 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%). For **Imagenet** [30], we use the off-the-shelf VGG16 model [31] as the black box. We select a sample of 100 images of the following classes French Bulldog, Scuba Diver, Corn, and Broccoli to use in the experiments. 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, the desired level of certainty is expressed as the width of the 95% credible interval.

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

Results for BayesLIME in Table 1 indicate that the true feature importances are close to ideal and indicate the estimates are well calibrated. While the estimates by BayesSHAP are somewhat less calibrated (true feature importances fall within our estimated 95% credible intervals about 89.2 to 98.4% of the time), they still are quite close to ideal. All in all, these results confirm that the credible intervals learned by our methods are well calibrated and therefore highly reliable in capturing the uncertainty of the feature importances. Lastly, though we set our priors to be uninformative in general, we also investigate how sensitive our uncertainty estimates are to hyperparameter choices in Figure 5 in the Appendix. We find that the explanation uncertainty becomes uncalibrated with strong priors. However, our explanations seem to be robust to hyperparameter choices in general.

**Correctness of Estimated Number of Perturbations** We assess whether our estimate of *perturbations-to-go* ( $G$ ; Section 3.2) is an accurate estimate of the *additional* number of perturbations needed to reach a desired level of feature importance certainty. We carry out this experiment on MNIST data for the digit “4” (additional datasets explored in Appendix C) and use  $S = 200$  as the initial number of perturbations to obtain a preliminary explanation and its associated uncertainty estimates. We then leverage these estimates to compute  $G$  for 6 different certainty levels. First, we observe significant differences in  $G$  estimates across instances (details in appendix C) i.e. number of perturbations needed to obtain a particular level of certainty varied significantly across instances—

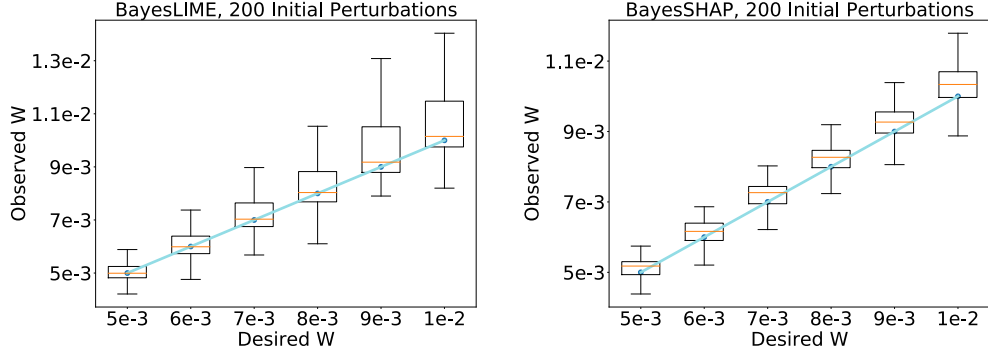


Figure 2: **Perturbations-to-go ( $G$ )**. We generate explanation with  $G$  perturbations, where  $G$  is computed using the *desired* credible interval width (x-axis), and compare desired levels to the *observed* credible interval width (y-axis) (blue line indicates ideal calibration). Results are averaged over 100 MNIST images of the digit “4”. We see that  $G$  provides a good approximation of the additional perturbations needed.

ranging from 200-5,000 for the lowest level of certainty to 200-20,000 for higher levels of certainty. Next, for each image and certainty level, we run our method for the estimated number of perturbations ( $G$ ) to determine if the observed estimates of uncertainty (observed credible interval width  $W$ ) match the desired levels of uncertainty (desired credible interval width  $W$ ). Results in Figure 2 show that the observed and desired levels of certainty are well calibrated, demonstrating that  $G$  estimates are reliable approximations of the additional number of perturbations needed.

**Efficiency of Focused Sampling** *Focused sampling* uses the *predictive variance* to strategically choose perturbations that will reduce uncertainty in order to be labeled by the black box (section 3.3). Here, we will evaluate the efficiency of the focused sampling procedure. First, we assess whether focused sampling converges (as measured by error uncertainty ( $P(\epsilon = 0)$ )) more efficiently than random sampling. To this end, we experiment with BayesLIME on Imagenet data for the “French bulldog” class to carry out this analysis. This setting replicates scenarios where LIME is applied to a computationally expensive black box model, making it highly desirable to limit the number of perturbations to reduce total running time. We run each sampling strategy for 2,000 perturbations and plot the number of model queries versus error uncertainty. During focused sampling, we set the batch size  $B$  to 50. The results in Figure 3 show that focused sampling results in faster convergence to reliable and high quality explanations; focused sampling stabilizes within a couple hundred model queries while random sampling takes over 1,000. Note, as the inefficiency of querying the black box model increases, the advantages of focused sampling decreasing total running time of the explanations will only become more pronounced. These results clearly demonstrate that focused sampling can significantly speed up the process of generating high quality local explanations. Additionally, in Appendix C, we also check if focused sampling causes any bias (due to sampling based on uncertainty estimates) that results in convergence to a different/wrong explanation, however our results clearly indicate that this is not the case.

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

$$\hat{L}(x_i) = \operatorname{argmax}_{x_j \in N_\epsilon(x_i)} \frac{\|\phi_i - \phi_j\|_2}{\|x_i - x_j\|_2} \quad (12)$$

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 parameters for  $x_i$  and  $x_j$ . Lower values indicate more stable explanations. We follow the setup outline by Alvarez-Melis and Jaakkola [18] and compute the local Lipschitz values, comparing both LIME & BayesLIME and SHAP & BayesSHAP across Compas, German Credit, MNIST digit “4”, and Imagenet “French Bulldog.” We perform the comparison using the default number of perturbations in



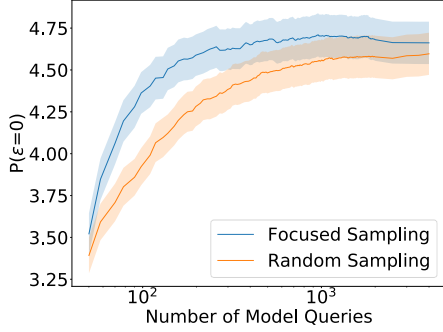


Figure 3: **Efficiency of focused sampling** for 100 Imagenet “French bulldog” images, with random sampling as a baseline. We provide mean and standard error. We assess the efficiency of focused sampling by comparing *error uncertainty* over model queries and show quicker convergence than random sampling.

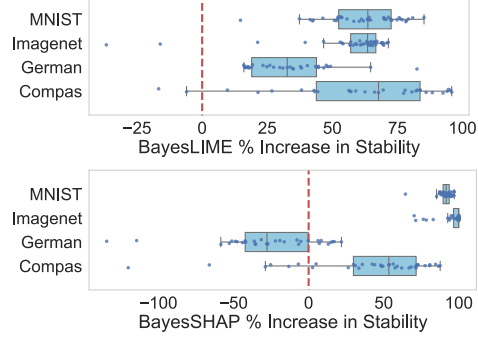


Figure 4: **Assessing the % increase in stability** of BayesLIME and BayesSHAP over LIME and SHAP respectively. Our Bayesian methods are significant more stable ( $\rho < 1e-2$  according to Wilcoxon signed-rank test) except for BayesSHAP on German Credit, where there is not a significant difference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

both LIME & SHAP, and use this same number in the respective Bayesian variants and set the batch size  $B$  to half this value. We use focused sampling for BayesLIME and BayesSHAP, and report the % increase in stability of these approaches over LIME and SHAP for 40 test points. The results given in Figure 4 show a clear improvement (on average 53%) in stability in all cases except German Credit for BayesSHAP. Further, we run a Wilcoxon signed-rank test and find our results are statistically significant in all cases ( $\rho < 1e-2$ ) except for BayesSHAP for German Credit, where there is not a significant difference between the methods ( $\rho > 0.05$ ). These results demonstrate BayesLIME and BayesSHAP are more stable than previous methods.

**User Study** We perform a user study with 31 subjects to compare BayesLIME and LIME explanations on MNIST. We evaluate the following: are explanations with low levels of uncertainty (i.e., most confident explanations) more meaningful to humans? To answer this question, we follow prior work and mask the most important features selected by BayesLIME and LIME [32, 4]. We ask users to guess the digit of the masked images. The better the explanation, the more difficult it should be for the users to get it right. Further, the choice to mask the important features is motivated by its success in prior work. We randomly select 15 correctly predicted test images, generate explanations by sweeping over a range of perturbation amounts  $[10^{-5}, \dots, 10^{3.5}]$  incremented by 0.5. We choose the *top* explanation for each image based on either fidelity (for LIME) or  $P(\epsilon = 0)$  (for BayesLIME). We sent the user study out to students and researchers with background in computer science. A screen shot of the task is shown in Figure 7 in the Appendix. We find that the explanations output by our methods focus on more informative parts of the image, since hiding them makes it difficult for humans to guess the digit. Users had an error rate of 25.7% for LIME, while it was 30.7% for BayesLIME, both with standard error 0.003 ( $\rho = 0.028$  through a one-tailed two sample t-test). This result indicates that our method BayesLIME and the associated measure of explanation uncertainty result in more high quality and reliable explanations compared to LIME and its associated fidelity metric.

## 5 Related Work

**Interpretability Methods** A variety of interpretability methods have been proposed. Some methods that are inherently interpretable include additive models [33, 34], decision lists and sets [35, 36], and instance-based explanations [37]. However, black-box models are often more flexible, accurate, and easier to use; thus, there has been a lot of interest in constructing post hoc explanations [38]. These include LIME [2] and SHAP [4, 39], which are among the most popular due to their broad applicability and code availability, but saliency maps [5–8], permutation feature importance [40], and partial dependency plots [41] also follow this paradigm. Other approaches to post hoc explanations focus on rule-based models [1, 3], counterfactuals [42, 43], and influence functions [9].

**Vulnerabilities of Post hoc Explanations** Recent work has shed light on the downsides of post hoc explanation techniques. These methods are often highly sensitive to small changes in inputs [14], are susceptible to manipulation [15, 16, 44, 45], and are not faithful to the underlying black boxes [46]. Perturbation-based explanation methods such as LIME and SHAP are subject to additional criticisms: results vary between runs of the algorithms [18–20, 47, 21], and hyperparameters used to select the perturbations can greatly influence the resulting explanation [20]. Prior work has attempted to tackle the problem of instability in perturbation-based explanations by averaging over several explanations [48, 19], however, this is computationally expensive. Other works related to creating more trustworthy explanations include development of sanity checks for explainers [49, 17, 50]. These techniques represent an important step towards improved usability, given experimental evidence that humans are often too eager to accept inaccurate machine explanations [51–54]. Recent works theoretically analyze the sources of non-robustness in black box explanations [55–57].

**Logical and Formal Reasoning** Additional related works have considered explaining classifiers through identifying a subset of features that are “sufficient” to explain a prediction [58–62]. Though these methods offer strong guarantees surrounding which features ensure a prediction is achieved, they are not model agnostic. Further, they do not define feature importances associated with the local explanations nor consider ways to improve locally weighted explanations, such as LIME and SHAP.

**Bayesian Methods in Explainable ML** Few recent works have adopted Bayesian formulations to explain black box models [63–65]. Guo et al. [63] introduce a Bayesian non-parametric approach to fit a *global* surrogate model. Their formulation seeks to fit a mixture of generalizable explanations across instances. Zhao et al. [64] study whether incorporating informative priors improves the stability of the resulting explanations. However, neither of these works focus on modeling the uncertainty of local explanations. Further, these approaches also do not tackle the critical problems of estimating key hyperparameters or improving efficiency of computing explanations.

## 6 Conclusion

We developed a Bayesian framework for generating local explanations along with their associated uncertainty. We instantiated this framework to obtain Bayesian versions of LIME and SHAP that output pointwise estimates of feature importances as well as their associated credible intervals. These intervals enabled us to infer the quality of the explanations and output explanations that satisfied user specified levels of uncertainty. We carried out theoretical analysis that leverages these uncertainty measures (credible intervals) to estimate the values of critical hyperparameters (e.g., the number of perturbations). We also proposed a novel sampling technique called focused sampling that leverages uncertainty estimates to determine how to sample perturbations for faster convergence.

While the Bayesian framework addresses several critical challenges (i.e., consistency, stability, modeling uncertainty) associated with LIME and SHAP, there are still certain aspects where it would exhibit the same shortcomings as LIME and SHAP [4, 66]. For instance, if the local decision surface of a given black box classifier is highly non-linear, our framework, which relies on local linear approximations, may not be able to capture this non-linear decision surface accurately. In addition, if the perturbation sampling procedures used in LIME and SHAP are used in BayesLIME and BayesSHAP, they will likely be vulnerable to the attacks proposed by Slack et al. [15]. In the future, it would be interesting to extend our framework to produce global explanations with uncertainty guarantees and explore how uncertainty quantification can help calibrate user trust in model explanations.

## 7 Acknowledgments

We would like to thank the anonymous reviewers for their insightful feedback. This work is supported in part by the NSF awards #IIS-2008461, #IIS-2008956, and #IIS-2040989, and research awards from the Harvard Data Science Institute, Amazon, Bayer, Google, and the HPI Research Center in Machine Learning and Data Science at UC Irvine. The views expressed are those of the authors and do not reflect the official policy or position of the funding agencies.

## References

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

- [18] David Alvarez-Melis and Tommi S. Jaakkola. On the robustness of interpretability methods. *ICML Workshop on Human Interpretability in Machine Learning*, 2018.
- [19] Eunjin Lee, David Braines, Mitchell Stiffler, Adam Hudler, and Daniel Harborne. Developing the sensitivity of lime for better machine learning explanation. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, page 1100610. International Society for Optics and Photonics, 2019.
- [20] Hui Fen Tan, Kuangyan Song, Madeilene Udell, Yiming Sun, and Yujia Zhang. “why should you trust my explanation?” understanding uncertainty in lime explanations. In *ICML Workshop on AI for Social Good*, 2019.
- [21] Jianbo Chen, Le Song, Martin J. Wainwright, and Michael I. Jordan. L-shapley and c-shapley: Efficient model interpretation for structured data. In *International Conference on Learning Representations*, 2019.
- [22] Andrew Moore. Locally weighted bayesian regression, January 1995.
- [23] Kacper Sokol, Alexander Hepburn, Raúl Santos-Rodríguez, and Peter A. Flach. blimey: Surrogate prediction explanations beyond lime. *NeurIPS HCML Workshop*, 2019.
- [24] Emily L. Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *NIPS*, 2014.
- [25] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. *BMVC 2014 - Proceedings of the British Machine Vision Conference 2014*, 05 2014.
- [26] Burr Settles. Active learning literature survey. 2010.
- [27] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias. In *ProPublica*, 2016.
- [28] Dheeru Dua and Casey Graff. Uci machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- [29] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- [30] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [32] Patrick Schwab and Walter Karlen. Cxplain: Causal explanations for model interpretation under uncertainty. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 10220–10230. Curran Associates, Inc., 2019. URL <http://papers.nips.cc/paper/9211-cxplain-causal-explanations-for-model-interpretation-under-uncertainty.pdf>.
- [33] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. Accurate intelligible models with pairwise interactions. In *KDD*, 2013.
- [34] Berk Ustun, Stefano Traca, and Cynthia Rudin. Supersparse linear integer models for interpretable classification. *arXiv preprint arXiv:1306.6677*, 2013.
- [35] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. Interpretable decision sets: A joint framework for description and prediction. In *Knowledge Discovery and Data mining (KDD)*, 2016.
- [36] Elaine Angelino, Nicholas Larus-Stone, Daniel Alabi, Margo Seltzer, and Cynthia Rudin. Learning certifiably optimal rule lists for categorical data. *arXiv preprint arXiv:1704.01701*, 2017.

- [37] Been Kim, Cynthia Rudin, and Julie Shah. The bayesian case model: A generative approach for case-based reasoning and prototype classification. *arXiv preprint arXiv:1503.01161*, 2015.
- [38] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-Agnostic Interpretability of Machine Learning. In *ICML Workshop on Human Interpretability in Machine Learning (WHI)*, June 2016.
- [39] Ian Covert and Su-In Lee. Improving kernelshap: Practical shapley value estimation using linear regression. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 3457–3465. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/covert21a.html>.
- [40] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [41] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [42] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable recourse in linear classification. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT\* ’19*, pages 10–19, 2019. ISBN 978-1-4503-6125-5.
- [43] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.
- [44] Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via adversarial model manipulation. In *Advances in Neural Information Processing Systems 32*, pages 2921–2932. 2019.
- [45] Junlin Wang, Jens Tuyls, Eric Wallace, and Sameer Singh. Gradient-based Analysis of NLP Models is Manipulable. In *Findings of the Association for Computational Linguistics: EMNLP (EMNLP Findings)*, page 247–258, 2020.
- [46] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206, 2019.
- [47] Muhammad Rehman Zafar and Naimul Mefraz Khan. Dlime: a deterministic local interpretable model-agnostic explanations approach for computer-aided diagnosis systems. *arXiv preprint arXiv:1906.10263*, 2019.
- [48] Chih-Kuan Yeh, Cheng-Yu Hsieh, Arun Suggala, David I Inouye, and Pradeep K Ravikumar. On the (in) fidelity and sensitivity of explanations. In *Advances in Neural Information Processing Systems*, pages 10965–10976, 2019.
- [49] Oana-Maria Camburu, Eleonora Giunchiglia, Jakob Foerster, Thomas Lukasiewicz, and Phil Blunsom. Can i trust the explainer? verifying post-hoc explanatory methods. *arXiv preprint arXiv:1910.02065*, 2019.
- [50] Mengjiao Yang and Been Kim. Bim: Towards quantitative evaluation of interpretability methods with ground truth. *arXiv:1907.09701*, 2019.
- [51] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. Interpreting interpretability: Understanding data scientists’ use of interpretability tools for machine learning. In *CHI*, April 2020.
- [52] Fred Hohman, Andrew Head, Rich Caruana, Robert DeLine, and Steven M Drucker. Gamut: A design probe to understand how data scientists understand machine learning models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2019.
- [53] Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman Vaughan, and Hanna Wallach. Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810*, 2018.

- [54] Himabindu Lakkaraju and Osbert Bastani. "how do i fool you?" manipulating user trust via misleading black box explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 79–85, 2020.
- [55] Damien Garreau and Ulrike von Luxburg. Looking deeper into lime. *arXiv preprint arXiv:2008.11092*, 2020.
- [56] Prasad Chalasani, Jiefeng Chen, Amrita Roy Chowdhury, Xi Wu, and Somesh Jha. Concise explanations of neural networks using adversarial training. In *International Conference on Machine Learning*, pages 1383–1391. PMLR, 2020.
- [57] Alexander Levine, Sahil Singla, and Soheil Feizi. Certifiably robust interpretation in deep learning. *arXiv preprint arXiv:1905.12105*, 2019.
- [58] Eric Wang, Pasha Khosravi, and Guy Van den Broeck. Probabilistic sufficient explanations. In *IJCAI*, 2021.
- [59] Andy Shih, Arthur Choi, and Adnan Darwiche. A symbolic approach to explaining bayesian network classifiers. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence, IJCAI'18*, page 5103–5111. AAAI Press, 2018. ISBN 9780999241127.
- [60] Weijia Shi, Andy Shih, Adnan Darwiche, and Arthur Choi. On tractable representations of binary neural networks. In Diego Calvanese, Esra Erdem, and Michael Thielscher, editors, *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning, KR 2020, Rhodes, Greece, September 12-18, 2020*, pages 882–892, 2020. doi: 10.24963/kr.2020/91. URL <https://doi.org/10.24963/kr.2020/91>.
- [61] Yacine Izza, Alexey Ignatiev, and Joao Marques-Silva. 2020.
- [62] Alexey Ignatiev, Nina Narodytska, and Joao Marques-Silva. Abduction-based explanations for machine learning models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33 (01):1511–1519, Jul. 2019. doi: 10.1609/aaai.v33i01.33011511. URL <https://ojs.aaai.org/index.php/AAAI/article/view/3964>.
- [63] Wenbo Guo, Sui Huang, Yunzhe Tao, Xinyu Xing, and Lin Lin. Explaining deep learning models – a bayesian non-parametric approach. In *Neural Information Processing Systems (NeurIPS)*. 2018.
- [64] Xingyu Zhao, Xiaowei Huang, Valentin Robu, and David Flynn. Baylime: Bayesian local interpretable model-agnostic explanations. *arXiv preprint arXiv:2012.03058*, 2020.
- [65] Kirill Bykov, Marina Höhne, Klaus-Robert Müller, Shinichi Nakajima, and Marius Kloft. How much can i trust you? – quantifying uncertainties in explaining neural networks. *arXiv*, 06 2020.
- [66] Sushant Agarwal, Shahin Jabbari, Chirag Agarwal, Sohini Upadhyay, Steven Wu, and Himabindu Lakkaraju. Towards the unification and robustness of perturbation and gradient based explanations. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 110–119. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/agarwal21c.html>.
- [67] Ludwig Fahrmeir, Thomas Kneib, and Stefan Lang. *Regression*. Statistik und ihre Anwendungen. Springer, 2007. ISBN 978-3-540-33932-8.
- [68] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.



## A Derivations

**Model Derivation** We write the joint posterior as

$$\phi, \sigma^2 | Y, \mathcal{Z} \propto \rho(Y|X, \beta, \sigma^2) \rho(\beta|\sigma^2) \rho(\sigma^2) \quad (13)$$

$$\begin{aligned} &\propto (\sigma^2)^{-N/2} \exp\left(-\frac{1}{2\sigma^2} (Y - \mathcal{Z}\phi)^T \text{diag}(\Pi_x(\mathcal{Z}))\right. \\ &\quad \left. (Y - \mathcal{Z}\phi)) (\sigma^2)^{-1} \exp\left(-\frac{1}{2\sigma^2} \phi^T \phi\right) (\sigma^2)^{-(1+\frac{n_0}{2})} \exp\left[\frac{-n_0\sigma_0^2}{2\sigma^2}\right] \end{aligned} \quad (14)$$

Letting  $\hat{\phi} = (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + I)^{-1} \mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) Y$ , we group terms in the exponentials according to  $\phi$ . The intermediate steps can be found in [67]. Suppressing dependence on  $Y$  and  $\mathcal{Z}$ , we can write down the conditional posterior of  $\phi$  as

$$\phi | \sigma^2 \propto \exp\left(\frac{1}{2}\sigma^{-2} [\phi - \hat{\phi}]^T (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + I) [\phi - \hat{\phi}]\right) \quad (15)$$

So, we can see that our estimates for the mean and variance of  $\rho(\phi|\sigma^2, Y, \mathcal{Z})$  are  $\hat{\phi}$  and  $\sigma^2 (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + I)^{-1}$ . Next, we derive the conditional posterior for  $\sigma^2$ . We identify the form of the scaled inverse- $\chi^2$  distribution in the joint posterior as in [22] and write

$$\sigma^2 | \hat{\phi} \sim \text{Inv-}\chi^2(N + n_0, \frac{n_0\sigma_0^2 + Ns^2}{n_0 + N}) \quad (16)$$

where  $s^2$  is defined as in equation 7.

**Derivation of equation 8** We establish the identity [22]:

$$\begin{aligned} \sigma^2 &\sim \text{Inv-}\chi^2(a, b) \text{ and } z | \sigma^2 \sim \mathcal{N}(\mu, \lambda\sigma^2) \\ &\iff z \sim t_{(v=a)}(\mu, \lambda b) \end{aligned} \quad (17)$$

We have,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ ,  $\sigma^2 \sim \text{Inv-}\chi^2(N + n_0, \frac{n_0\sigma_0^2 + Ns^2}{n_0 + N})$ . Then, it's the case that  $\epsilon \sim t_{(v=N+n_0)}(0, \frac{n_0\sigma_0^2 + Ns^2}{n_0 + N})$ .

**Derivation of Posterior Predictive** Note, this derivation takes the priors to be set as in BayesLIME or BayesSHAP, namely, with values close to zero. We apply the identity from equation 17 to derive this posterior. We have  $\hat{y} \sim \hat{\phi}^T z + \epsilon$  for some  $z$ . Thus,  $\hat{y} \sim \mathcal{N}(\hat{\phi}^T z, z^T V_\phi z \sigma^2) + \mathcal{N}(0, \sigma^2)$ , where  $\sigma^2 \sim \text{Inv-}\chi^2(N, s^2)$ . So, we have  $\hat{y} \sim t_{(v=N)}(\hat{\phi}^T z, (z^T V_\phi z + 1)s^2)$ .

## B Proof of Theorems

In these derivations, the perturbation matrices  $\mathcal{Z}$  have elements  $\mathcal{Z}_{ij} \in \{0, 1\}$  where each  $\mathcal{Z}_{ij} \sim \text{Bernoulli}(0.5)$ . Note, in these proofs, we take the priors to be set as in BayesLIME and BayesSHAP, i.e., they have hyperparameter values close to 0.

### B.1 Proof of Theorem 3.3

Note that we use  $N$  to denote the *total* perturbations while  $S$  denotes the perturbations collected *so far*. We use three assumptions stated as follows. First,  $\frac{\pi N}{2}$  is sufficiently large such that  $\frac{\pi N}{2} + 1$  is equivalent to  $\frac{\pi N}{2}$ . Second,  $N$  is sufficiently large such that  $N + 1$  is equivalent to  $N$  and  $\frac{N}{N-2}$  is equivalent to 1. Third, the product of  $\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z}$  within  $V_\phi$  can be taken at its expected value. First, we state the marginal distribution over feature importance  $\phi_i$  where  $i$  is an arbitrary feature importance  $i \in d$ . This given as

$$\phi_i | \mathcal{Z}, Y \sim t_{v=N}(\hat{\phi}_i, V_{\phi_{ii}} s^2) \quad (18)$$

where  $V_\phi = (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) \mathcal{Z} + I)^{-1}$ . Recall each  $\mathcal{Z}_{ij}$  is given  $\sim \text{Bern}(0.5)$  we use the third assumption to write  $V_\phi$  is  $\frac{\pi N}{2} + 1$  for the on diagonal elements and  $\frac{\pi N}{4}$  for the off diagonal elements.

We can see this is the case considering that each element in  $\mathcal{Z}$  is a  $\text{Bern}(.5)$  draw. We drop the  $1's$  due to the first assumption.

Let  $k = \frac{\bar{\pi}N}{2}$ . It follows directly from Sherman Morrison that the  $i$ -th and  $j$ -th entries of  $V_\phi$  are given as

$$(V_\phi)_{ij} = \begin{cases} \frac{2}{k} - \frac{2}{k(N+1)} & i = j \\ -\frac{2}{k(N+1)} & i \neq j \end{cases} \quad (V_\phi)_{ii} = \frac{4}{\bar{\pi}(N+1)} \quad (19)$$

We see that the diagonals are the same. Thus, we take the  $PTG$  estimate in terms of a single marginal  $\phi_i$ . Substituting in the  $s^2$  estimate  $s_S^2$  and using the second assumption, we write the variance of marginal  $\phi_i$  as

$$\text{Var}(\phi_i) = \frac{4s_S^2}{\bar{\pi}(N+1)} \frac{N}{N-2} \quad (20)$$

$$= \frac{4s_S^2}{\bar{\pi} \times N} = \frac{4s_S^2}{\bar{\pi} \times \text{Var}(\phi_i)} \quad (21)$$

Because feature importance uncertainty is in the form of a credible interval, we use the normal approximation of  $\text{Var}(\phi_i)$  and write

$$N = \frac{4s_S^2}{\bar{\pi} \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} \quad (22)$$

where  $W$  is the desired width,  $\alpha$  is the desired confidence level, and  $\Phi^{-1}(\alpha)$  is the two-tailed inverse normal CDF. Finally, we subtract the initial  $S$  samples.  $\square$

## B.2 Proposition 3.2

Before providing a proof for proposition 3.2, we note to readers that the claims are related to well known results in bayesian inference (e.g. similar results are proved in [68]). We provide the proofs here to lend formal clarity to the properties of our explanations.

**Convergence of  $\text{Var}(\phi)$**  Recall the posterior distribution of  $\phi$  given in equation 5. In equation 19, we see the on and off-diagonal elements of  $V_\phi$  are given as  $\frac{4}{\bar{\pi}(N+1)}$  and  $-\frac{4}{\bar{\pi}N(N+1)}$  respectively (here replacing  $S$  with  $N$  to stay consistent with equation 5). Because we have  $N \rightarrow \infty$ , these values define  $V_\phi$  due to the law of large numbers. Thus, as  $N \rightarrow \infty$ ,  $V_\phi$  goes to the null matrix and so does the uncertainty over  $\phi$ .

**Consistency of  $\hat{\phi}$**  Recall the mean of  $\phi$ , denoted  $\hat{\phi}$  given in equation 6. To establish consistency, we must show that  $\hat{\phi}$  converges in probability to the true  $\phi$  as  $N \rightarrow \infty$ . To avoid confusing true  $\phi$  with the distribution over  $\phi$ , we denote the true  $\phi$  as  $\phi^*$ . Thus, we must show  $\hat{\phi} \rightarrow_p \phi^*$  as  $N \rightarrow \infty$ . We write

$$\hat{\phi} = (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))\mathcal{Z} + I)^{-1} \mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))Y \quad (23)$$

$$= (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))\mathcal{Z} + I)^{-1} \mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))(\mathcal{Z}\phi^* + \epsilon) \quad (24)$$

Considering mean of  $\epsilon$  is 0 and using law of large numbers,

$$= (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))\mathcal{Z} + I)^{-1} \mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))\mathcal{Z}\phi^* = \phi^* \quad (25)$$

**Convergence of  $\text{Var}(\epsilon)$**  Assume we have  $N \rightarrow \infty$  so  $\hat{\phi}$  converges to  $\phi^*$ . The uncertainty over the error term is given as the variance of the distribution in equation 8. The variance of this generalized student's t distribution is given as converges to  $s^2$  for large  $N$ . Recalling its definition,  $s^2$  reduces to the local error of the model as  $N \rightarrow \infty$ . which is equivalent to the squared bias of the local model.

## C Detailed Results

In this appendix, we provide extended experimental results.

### C.1 Explanation Uncertainty Hyperparameter Sensitivity

In the main paper, we assume the priors are set to be uninformative. Though this is the advised configuration for BayesLIME and BayesSHAP because prior information about the local surface is not likely available, we assess the calibration sensitivity of BayesLIME to different choices in the hyperparameters. In figure 5, we perform a grid search over the uncertainty hyperparameters  $n_0$  and  $\sigma_0^2$  for the MNIST digit “4” class. We find the explanation uncertainty is robust to the choice of hyperparameters.

$n_0 \backslash \sigma_0^2$	$1e-5$	$1e-1$	1	10	100
$1e-5$	95.7	95.7	96.4	96.6	96.6
$1e-1$	96.6	96.6	96.9	98.9	100.0
1	96.5	96.9	98.6	100.0	100.0
10	94.2	98.2	100.0	100.0	100.0
100	72.2	99.0	100.0	100.0	100.0

Figure 5: BayesLIME calibration sensitivity to the choice of hyperparameters. Closer to 95.0 is better. These results indicate BayesLIME calibration is not very sensitive to choices in the hyperparameter values.

### C.2 PTG Estimate Results

In the main paper, we provided PTG results for BayesLIME on MNIST. In this appendix, we show the number of perturbations estimated by PTG and additional PTG results on the Imagenet “French bulldog” class.

**Number of Perturbations Estimated by PTG** In section 4, we assessed if  $G$  produces good estimates of the number of additional samples needed to reach the desired level of feature importance certainty. In figure 6, we show the desired level of certainty (desired width of credible interval  $W$ ) versus the actual  $G$  estimate (i.e. the estimated number of perturbations) for figure 2 in the main paper. We see the estimated number of perturbations is highly variable depending on desired  $W$ .

**Further PTG Estimate Results** We provide results for the PTG estimate on Imagenet in Figure 8. We limit the range of uncertainty values compared to MNIST because the Imagenet data is more complex and consequently the required number of perturbations becomes very high. These results further indicate the effectiveness of the PTG estimate.

### C.3 User Study

**Participate Consent** We sent out an email to students and researchers with a background in computer science inviting them to take our user study. At the beginning of the user study, we stated that no

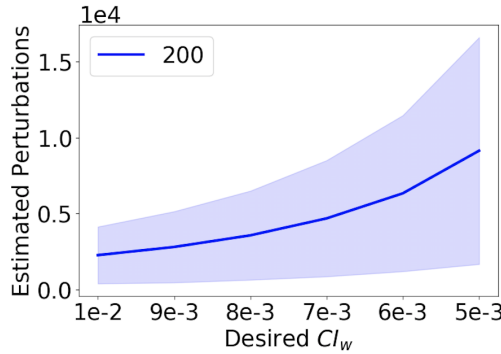


Figure 6: Desired  $CI_w$  versus the actual number of perturbations estimated by PTG in figure 2 of the main paper. We plot mean and standard deviation of  $G$ .

personal information would be asked during the study, their answers would be used in a research project, and whether they consented to take the study.

**Image of interface** We give an example screen shot from the user study in Figure 7.



Figure 7: Screen shot from user study (correct answer 4).

**BayesLIME Toy Example** To show how our bayesian methods capture the uncertainty of local explanations, we provide an illustrative example in Figure 9. Rerunning LIME explanations on a toy decision surface (blue lines in the figure), we see LIME has high variance and produces many different explanations. This behavior is particularly severe in the nonlinear surfaces. With a single explanation, BayesLIME captures the uncertainty associated with generating local explanations (black lines in the figure).

#### C.4 Focused Sampling Results

In this appendix, we provide additional focused sampling results. We include a comparison of focused sampling to random sampling in terms of wall clock time. We also provide results demonstrating the focused sampling procedure is not biased.

**Wall Clock Time of Focused Sampling** In figure 10, we plot wall clock time versus  $P(\epsilon = 0)$ . This experiment is analogous to figure 3 in the main paper, but here we use time instead of number of

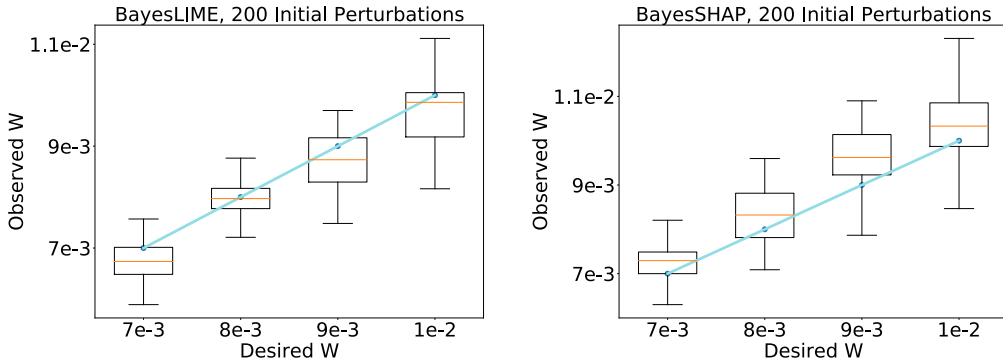
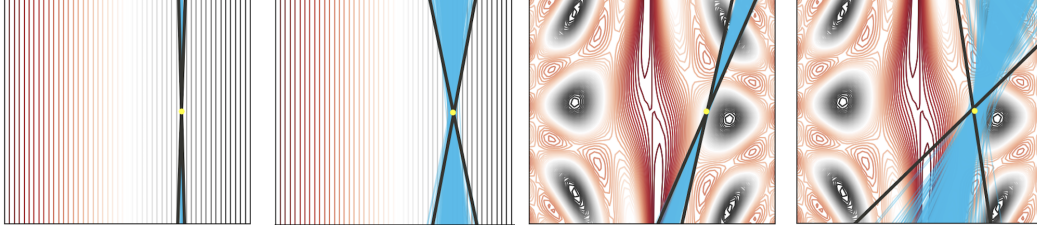


Figure 8: Imagenet PTG results for BayesLIME & BayesSHAP. The blue line indicates ideal calibration. These results indicate the PTG estimate is well calibrated for BayesLIME and BayesSHAP on Imagenet, demonstrating the efficacy of the estimate.



(a) linear, many samples (b) linear, fewer samples (c) nonlinear, many samples (d) nonlinear, fewer samples

Figure 9: Rerunning LIME local explanations 1000 times and BayesLIME *once* for linear and non-linear toy surfaces using few (25) and many (250) perturbations. The linear surface is given as  $p(y) \propto x_1$  and the non linear surface is defined as  $p(y) \propto \sin(x_1/2) * 10 + \cos(10 + (x_1 * x_2)/2) * \cos(x_1)$ . We plot each run of LIME in blue and the BayesLIME 95% credible region of the feature importance  $\phi$  in black. We see that LIME variance is higher with fewer samples and a less linear surface. BayesLIME captures the relative difficulty of explaining each surface through the width the credible region. For instance, BayesLIME is most uncertain in the nonlinear, few samples case because this surface is the most difficult to explain.

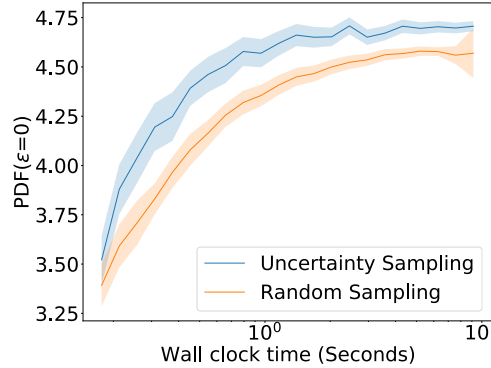


Figure 10: Wall clock time needed to converge to a high quality explanation by BayesLIME (analogous to figure 3 in the paper). We use both random sampling and focused sampling over 100 Imagenet images. We provide the mean and standard error for binned estimates of these values. This result demonstrates that focused sampling leads to improved convergence over random sampling in terms of wall clock time.

model queries on the x-axis. We see that uncertainty sampling is more time efficient than random sampling for BayesLIME.

**Bias of Focused Sampling** In the main text, we saw that focused sampling converges faster than random sampling. However, it is possible that focused sampling introduces bias into the process due to sampling based on uncertainty estimates, leading to convergence to a different/wrong explanation. To assess whether this occurs in practice, we evaluate the convergence of both focused sampling and random sampling to the “true” explanation on Imagenet (computed with the number of perturbations  $N = 10,000$  using random sampling). To measure convergence, we compare the  $L_1$  distance of the explanation with the ground truth explanation. The result provided in Figure 11 demonstrates that focused sampling converges to the ground truth explanation with significantly fewer model queries than random sampling. Focused sampling reaches a  $L_1$  distance of 0.1 at 300 queries while it takes upwards of 450 queries for random sampling, indicating improved query efficiency of 30 – 40%. Lastly, as the number of model queries increases ( $\sim 1000$ ), we observe an  $L_1$  distance of around 0.06 which is extremely small and the explanations are practically the same as the ground truth. Overall, these results show that focused sampling does not suffer from biases in practice and further demonstrate that focused sampling can lead to significant speedups.

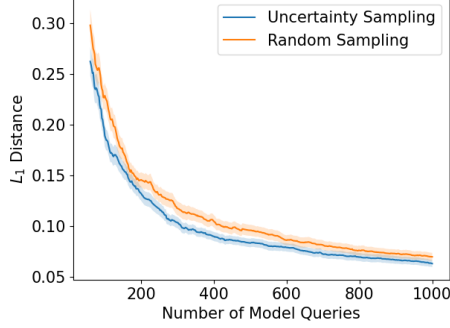


Figure 11: Convergence of both focused sampling and random sampling to the ground truth explanation. We see that uncertainty sampling converges more quickly to the ground truth than random sampling, demonstrating that there is minimal bias in the focused sampling procedure and focused sampling converges more efficiently.

### C.5 Benchmarking

We also benchmark the efficiency of BayesLIME and BayesSHAP against Guo et al. [63], a related Bayesian explanation method that uses a Bayesian non parametric mixture regression and MCMC for parameter inference. Fixing their mixture regression to a single component results in a similar model to ours and thus is a useful point of comparison. To explain a single instance on ImageNet using VGG16, their approach takes 139.2 seconds, while BayesLIME and BayesSHAP take 20.3 seconds and 21.1 seconds respectively, under the same conditions, demonstrating that the closed form solution is very efficient.

## D Explaining a Ground Truth Function

We consider a synthetic experiment in which we observe an underlying ground truth function and verify that lower values of feature importance uncertainty indicate higher proximity between the feature importance estimates and the underlying ground truth function. To this end, we constructed a piecewise linear function of two variables, where each quadrant in the  $x,y$ -plane corresponds to a different linear model. We consider the regression coefficients of the quadrant as the ground truth explanation. The piecewise function is given as:

$$\begin{aligned}
 f(x, y) = & 0.3x + 0.2y \text{ if } x > 0, y > 0 \\
 & 0.2x - 0.1y \text{ if } x > 0, y \leq 0 \\
 & -x - 0.05y \text{ if } x \leq 0, y \leq 0 \\
 & -.8x + 0.2y \text{ if } x \leq 0, y > 0
 \end{aligned} \tag{26}$$

We plot the  $\ell_1$  distance between the BayesLIME feature importance mean and ground truth explanation versus the maximum credible interval width of the BayesLIME explanation. The results given in Figure 12 indicate that tighter credible intervals lead to explanations that are closer to the ground truth, demonstrating that the feature importance uncertainties are meaningful in regards to a ground truth function.

## E Compute Used

In this work, we ran all experiments on a single NVIDIA 2080TI & a single NVIDIA Titan RTX GPU.

## F Dataset licenses

German Credit is in the public domain, COMPAS uses the MIT license, MNIST uses the Creative Commons Attribution-Share Alike 3.0 license, and Imagenet does not hold copyright of images.



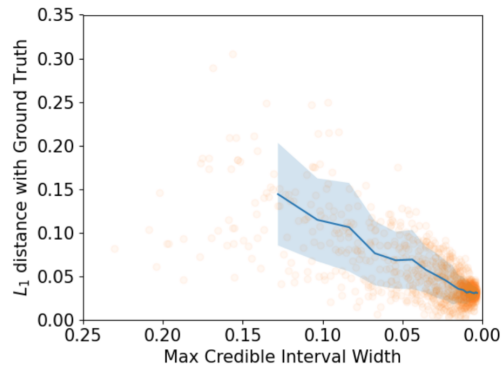


Figure 12: Assessing whether tighter credible intervals lead to convergence with ground truth, on an example where the ground truth feature importances are known. Here, we plot The  $\ell_1$  distance between the feature importances means for BayesLIME and ground truth explanation versus the maximum credible interval width across the explanation.