

## **Paper Review**

# **Reliable Post hoc Explanations: Modeling Uncertainty in Explainability**

**NeurIPS 2021**

**Dylan Slack<sup>1</sup> , Sophie Hilgard<sup>2</sup>, Sameer Singh<sup>1</sup> , Himabindu Lakkaraju<sup>2</sup>**

**<sup>1</sup>UC Irvine <sup>2</sup>Harvard University**

Min-Seok Yang

Natural Language Processing Lab

Department of Artificial Intelligence, Kyung Hee University

# Background

---

- **Explain machine learning model in deployment**
  - Model deployment in domains such as healthcare, medicine, law, finance
    - Important to ensure that decision makers have a clear understanding of the behavior of these models
  - Explain complex black box models by constructing interpretable local approximations
    - Post-hoc explanations construct interpretable local approximations
    - Lime, SHAP, MAPLE, Anchors

## Post-hoc interpretability

- Interpretability is achieved by applying methods that analyze the model after training

## Local Method

- Model-agnostic methods which focus on explaining individual prediction

# Background

---

- **Existing local explanation methods**
  - Lime, SHAP, MAPLE, Anchors
  - Unstablility
    - Negligibly small perturbations to an instance can result in substantially different explanation
  - Inconsistency
    - Multiple runs on the same input instance with the same parameter settings may result in vastly different explanations

## Part 1. Background

- **Reliable metrics to ascertain the quality of the explanations**
  - Explanation fidelity rely heavily on the implementation details of explanation method
    - 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
    - Typically require a large number of black box model queries to construct local approximations

### Explanation produced by LIME

$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

Minimize  $\mathcal{L}(f, g, \pi_x)$  while having  $\Omega(g)$  be low enough to be interpretable by humans

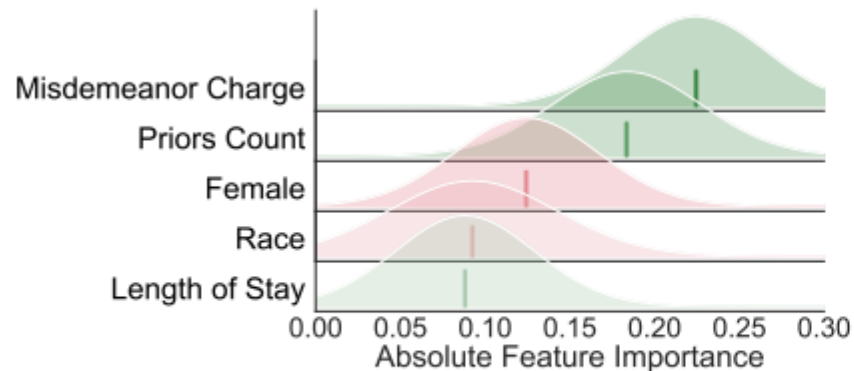
### Formulation

- Explanation family  $G$
- Fidelity function  $\mathcal{L}$
- Complexity measure  $\Omega$

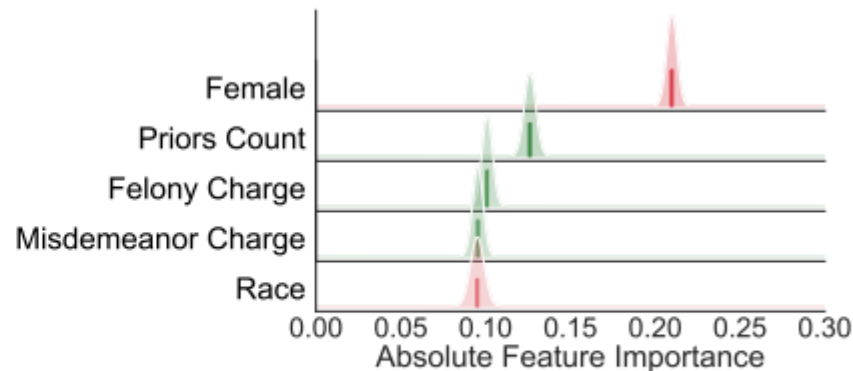
## Part 2. Introduction

- **Bayesian framework for generating local explanations with uncertainty**

- Generate consistent, stable, and reliable explanations with guarantees in a computationally efficient manner
- 1. Bayeslime & BayesSHAP: Bayesian versions of LIME and KernelSHAP



(a) Explanation computed with 100 perturbations



(b) Explanation with 2000 perturbations

### Example Explanation

- Dataset: COMPAS
- Vertical Lines: LIME
- Shaded Region: BayesLIME
- Red: Negative Effect
- Green: Positive Effect

**LIME:** Contradictory feature importance for different number of perturbations

**BayesLIME:** Provide more context (i.e., tighter uncertainty interval indicates importance)

- **Bayesian framework for generating local explanations with uncertainty**
  - Generate consistent, stable, and reliable explanations with guarantees in a computationally efficient manner
  - 1. Bayeslime & BayesSHAP: Bayesian versions of LIME and KernelSHAP
  - 2. Closed form expressions for the posteriors of the explanations
    - Eliminate the need for any additional computational complexity
  - 3. Credible intervals produced by our framework
    - Make concrete inferences about the quality of the resulting explanations
    - Produce explanations that satisfy user specified levels of uncertainty (e.g., 95% confidence level)

- **Notation**

- $f : \mathbb{R}^d \rightarrow [0, 1]$

- Black box classifier that takes a data point  $x$  with  $d$  features
    - Returns the probability that  $x$  belongs to a certain class

- $\phi \in \mathbb{R}^d$

- 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

- $\phi$  captures the coefficients of a linear model

- Let  $\mathcal{Z}$  be a set of  $N$  randomly sampled instances (perturbations) around  $x$

- Proximity between  $x$  and any  $z \in \mathcal{Z}$  is given by  $\pi_x(z) \in \mathbb{R}$

- Vector of these distances over  $N$  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 & KernelSHAP**

- Model-agnostic local explanation approaches that explain predictions of a classifier  $f$  by learning a linear model locally  $\phi$  around each prediction (i.e.  $y \sim \phi^T z$ )
- Objective function

- Explanation that approximates the behavior of the black box accurately in the vicinity (neighborhood) of  $x$

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

- Closed form solution from objective function

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

- LIME

- Chosen  $\pi_x(z)$  heuristically: Cosine or  $l_2$  distance

- KernelSHAP

- Game theoretic principles to compute  $\pi_x(z)$ , guaranteeing that explanations satisfy certain properties



- **Constructing Bayesian Local Explanations**

- 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  $\epsilon$
- Weights of linear combination  $\phi$  capture feature importances & constitute explanation
- $\epsilon$  captures the error that arises due to mismatch between explanation  $\phi$  & local decision surface of the black box model  $f$

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

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

- **Constructing Bayesian Local Explanations**

- Error term is modeled as a Gaussian whose variance relies on proximity function  $\pi_x(z)$
- Proximity function  $\pi_x(z)$ 
  - Perturbations closer to the data point  $x$  are modeled accurately
  - Allows more room for error in case of perturbations that are farther away
  - Cosine or  $l_2$  distance or game theoretic principles similar to LIME & KernelSHAP
- Distributions on error  $\epsilon$  and feature importance  $\phi$  both consider the parameter  $\sigma^2$
- 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
  - Expect to be more confident about the feature importances

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

- **Constructing Bayesian Local Explanations**

- The weighted least squares formulation of LIME and KernelSHAP

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

- Corresponds to the Bayesian version of that of LIME and KernelSHAP with additional terms to model uncertainty

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

- Feature importance uncertainty
  - The uncertainty associated with the feature importances  $\phi$
- 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.

- **Constructing Bayesian Local Explanations**

- Inference process involves estimating the values of two key parameters:  $\phi$  and  $\sigma^2$ 
  - Compute the local explanation as well as the uncertainties associated with feature importances  $\phi$  and the error term  $\epsilon$
- Posterior distributions on  $\phi$  and  $\sigma^2$  are normal and scaled **Inv- $\chi^2$** , respectively, due to the corresponding conjugate priors

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

$$\phi | \sigma^2, \mathcal{Z}, Y \sim \text{Normal}(\hat{\phi}, V_{\phi} \sigma^2)$$

$$\hat{\phi} = V_{\phi} (\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z})) Y)$$

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

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

- **Constructing Bayesian Local Explanations**

- Estimate of the posterior mean feature importances of LIME and KernelSHAP

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

- Our Estimate of the posterior mean feature importances

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

- If use the same proximity function  $\pi_x(z)$  in our framework as in LIME & KernelSHAP, the posterior mean of the feature importance  $\hat{\phi}$  output by our framework will be equivalent to the feature importances output by LIME & KernelSHAP

- **Constructing Bayesian Local Explanations**

- Feature Importance Uncertainty

- Compute the posterior mean of local feature importances  $\hat{\phi}$  using closed form expression
    - Estimate credible interval (measure of uncertainty) around the mean feature importances by repeatedly sampling from the posterior distribution

**Closed form expression**

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

- **Constructing Bayesian Local Explanations**

- Error Uncertainty

- Error term  $\epsilon$  can serve as a proxy for explanation quality
    - Captures the mismatch between the constructed explanation and the local decision surface of the underlying black box

Equation  $s^2$

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

Student's t distribution

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

- Evaluate the probability density function of the above posterior at 0, i.e.,  $P(\epsilon = 0)$
    - Substituting the value of  $s^2$  computed using equation  $s^2$  into the Student's t distribution
    - Perfectly captures the local decision surface underlying the black box
    - Operation in constant time, adding minimal overhead to non-Bayesian LIME & SHAP

- **BayesLIME**

- Obtain the Bayesian version of LIME by setting the proximity function
- $D$  Distance metric (e.g. cosine or  $l_2$  distance) ,  $n_0$  &  $\sigma_0^2$  to small values (  $10^{-6}$  )
- Prior is uninformative
- Compute feature importance uncertainty & error uncertainty for LIME's feature importances

**Proximity function**

$$\pi_x(z) = \exp(-D(x, z)^2 / \sigma^2)$$



- **BayesSHAP**

- Obtain the Bayesian version of KernelSHAP by setting uninformative prior on  $\sigma^2, \pi_x(z)$
- SHAP method views the problem of constructing a local linear model as estimating the Shapley values corresponding to each of the features
- Shapley values represent the contribution of each of the features to the black box prediction
- The measures of uncertainty output by our method BayesSHAP capture the reliability of the estimated variable contributions

**Proximity function**

$$\pi_x(z) = \frac{d-1}{(d \text{ choose } |z|)|z|(d-|z|)}$$

$|z|$  The number of the variables in the variable combination represented by the data point  $z$

- **BayesLIME & BayesSHAP**
  - Encourage BayesLIME and BayesSHAP explanations to be sparse
    - Use dimensionality reduction or feature selection techniques as used by LIME and SHAP to obtain the top K features
    - Construct our explanations using the data corresponding to these top K features

- **Estimating the Number of Perturbations**

- Major drawbacks of approaches such as LIME and KernelSHAP
  - Do not provide any guidance on how to choose the number of perturbations, a key factor in obtaining reliable explanations in an efficient manner
- Leverage the uncertainty estimates output by our framework to compute perturbations-to-go  $G$
- Perturbations-to-go  $G$  : Estimate of how many more perturbations are required to obtain explanations that satisfy a desired level of certainty
  - 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  $W = 0.1$ .

- **Estimating the Number of Perturbations**
  - Estimate  $G$  for the local explanation of a data point  $\mathcal{X}$ 
    - Generate  $S$  perturbations around  $\mathcal{X}$  (where  $S$  is small and chosen by the user)
    - Fit a local linear model using our method
    - Provides initial estimates of various parameters shown in equations

#### Equations

$$\hat{\phi} = V_{\phi}(\mathcal{Z}^T \text{diag}(\Pi_x(\mathcal{Z}))Y)$$

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

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

## • Estimating the Number of Perturbations

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

$$\left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2 = \text{Var}(\phi_i) = \frac{4s_S^2}{\bar{\pi}_S \times (G + S)} \quad \Rightarrow \quad \text{Perturbations-to-go } G \quad G(W, \alpha, x) = \frac{4s_S^2}{\bar{\pi}_S \times \left[ \frac{W}{\Phi^{-1}(\alpha)} \right]^2} - S$$

- Average proximity  $\pi_x(z)$  for the  $S$  perturbations  $\bar{\pi}_S$
- Empirical sum of squared errors between black box & local linear model predictions  $s_S^2$
- The two-tailed inverse normal CDF at confidence level  $\alpha$

- **Focused Sampling of Perturbations**

- If Perturbations-to-go  $G$  is large, Need to reduce this cost
  - Focused sampling which leverages uncertainty estimates to query the black box in a more targeted fashion (instead of querying randomly)
- Inspired by active learning, focused sampling strategically prioritizes perturbations whose predictions the explanation is most uncertain about
  - Query the black box only for the predictions of the most informative perturbations
  - Learn an accurate explanation with far fewer queries to the black box
- Determine how uncertain our explanation  $\phi$  is about the black box label for any given instance Query  $z$ 
  - Compute the posterior predictive distribution for  $z$   
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

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

- **Focused Sampling of Perturbations**

- The variance of this three parameter student's t distribution

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

- Refer to this variance as the predictive variance  $\text{var}(\hat{y}(z))$
- Captures how uncertain our explanation  $\phi$  is about the black box prediction

---

**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
```

---

# Experiment

---

- **Process**

- Evaluate proposed framework by first analyzing the quality of our uncertainty estimates
  - i.e., feature importance uncertainty and error uncertainty
- Assess our estimates of required perturbations  $G$
- Evaluate the computational efficiency of focused sampling
- Describe user study with 31 subjects to assess the informativeness of the explanations output by our framework



# Experiment

---

- **Dataset**

- COMPAS
  - Criminal history, jail and prison time, and demographic attributes of 6172 defendants
- German Credit
  - Financial and demographic information for 1000 loan applications, each labeled as a "good" or "bad" customer
- MNIST
  - Handwritten digits dataset
- Imagenet
  - Select a sample of 100 images of classes French Bulldog, Scuba Diver, Corn, and Broccoli

- **Model**

- Random forest classifier (sklearn implementation with 100 estimators) as black box models for COMPAS, German Credit
- 2-layer CNN to predict the digits for MNIST
- The off-the-shelf VGG16 model as the black box

# Experiment

---

- **Baseline**

- Generating explanations
  - LIME and KernelSHAP with default settings
- Images
  - Construct super pixels as described in LIME
  - Use them as features (number of super pixels is fixed to 20 per image)
- Parameter
  - The desired level of certainty is expressed as the width of the 95% credible interval

## Part 4. Experiment

- **Quality of Uncertainty Estimates**

- Well calibrated
- Highly reliable in capturing the uncertainty of the feature importances

95% credible intervals with 100 perturbations include their true values (estimated on 10,000 perturbations)

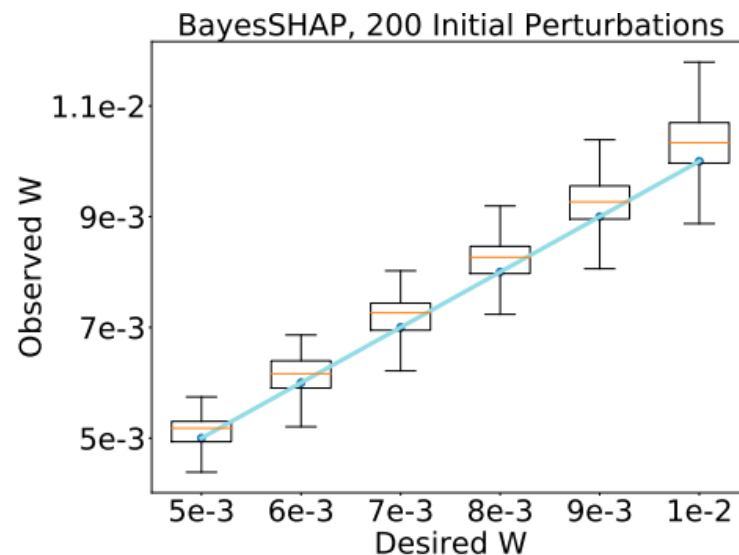
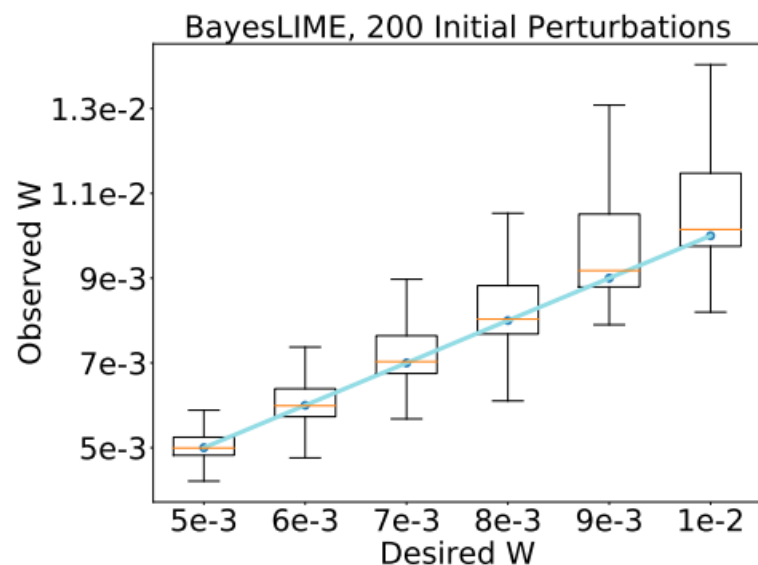
	BayesLIME	BayesSHAP		BayesLIME	BayesSHAP
TABULAR DATASETS			MNIST		
COMPAS	95.5	87.9	Digit 1	95.8	98.4
German Credit	96.9	89.6	Digit 2	95.8	97.4
IMAGENET			Digit 3	95.2	96.3
Corn	94.6	91.8	Digit 4	97.2	90.1
Broccoli	91.4	89.2	Digit 5	95.2	95.6
French Bulldog	94.8	89.9	Digit 6	96.7	96.8
Scuba Diver	92.4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

## Part 4. Experiment

### • Correctness of Estimated Number of Perturbations

- Leverage these estimates to compute  $G$  for 6 different certainty levels
- Perturbations, where  $G$  is computed using the desired credible interval width (x-axis),
- Compare desired levels to the observed credible interval width (y-axis)
- Blue line indicates ideal calibration
- Provides a good approximation of the additional perturbations needed

Perturbations-to-go (  $G$  ) is averaged over 100 MNIST images of the digit "4"

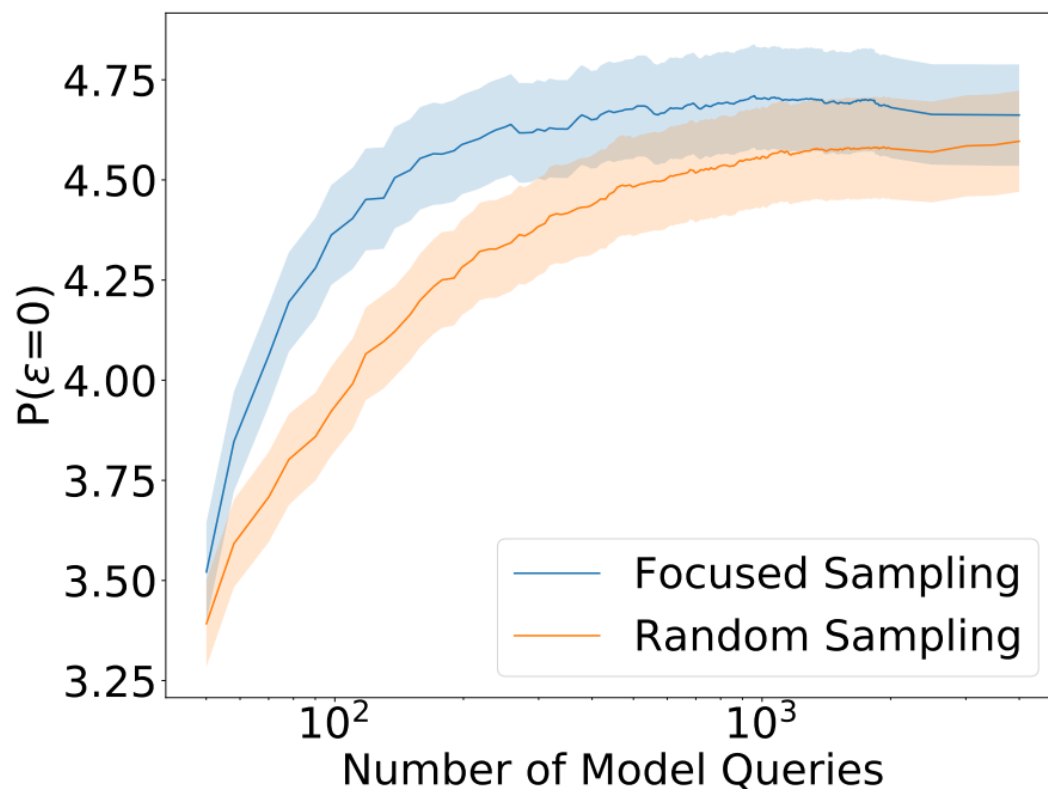


$$\hat{L}(x_i) = \underset{x_j \in N_\epsilon(x_i)}{\operatorname{argmax}} \frac{\|\phi_i - \phi_j\|_2}{\|x_i - x_j\|_2}$$
$$N_\epsilon(x_i) \quad x_i \quad \phi_i \quad \phi_j$$

## Part 4. Experiment

- **Efficiency of Focused Sampling**

- Results in faster convergence to reliable and high quality explanations
- Stabilizes within a couple hundred model queries while random sampling takes over 1,000



## Part 4. Experiment

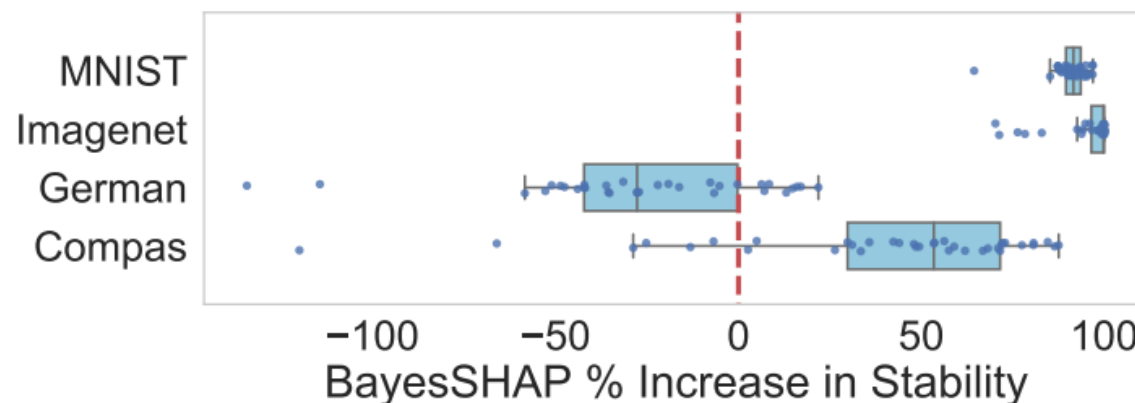
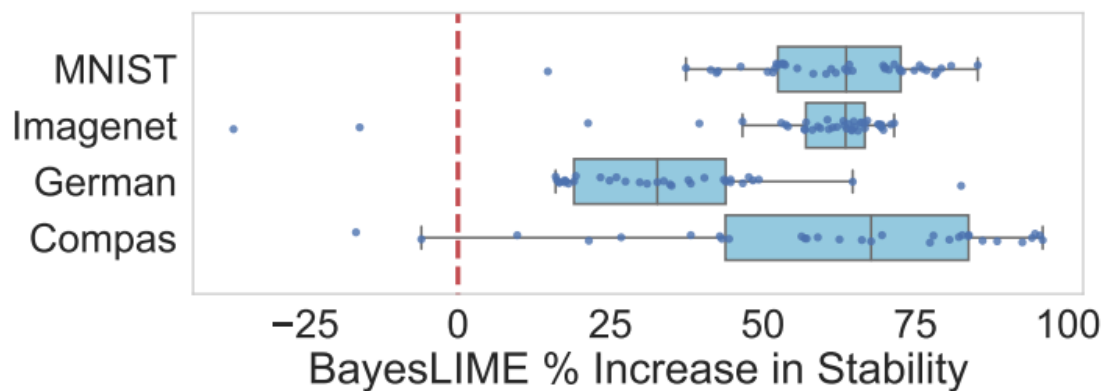
- **Stability of BayesLIME & BayesSHAP**

- Use the local Lipschitz metric for explanation stability

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

- Clear improvement (on average 53%) in stability in all cases except German Credit for BayesSHAP

**Assessing the % increase in stability of BayesLIME and BayesSHAP over LIME and SHAP respectively**



## Part 5. Conclusion

- **Bayesian framework as solution of previous local explanations**

Problem	Solution
Difficult to set hyperparameters	PTG(Perturbations-to-go)
Unclear when you have a good explanation	Credible Intervals
Unstable, re-reruns lead to different explanations	Credible Intervals
Often naive sampling : Focused sampling	Focused sampling

# Thank you

---

2023. 05. 15.

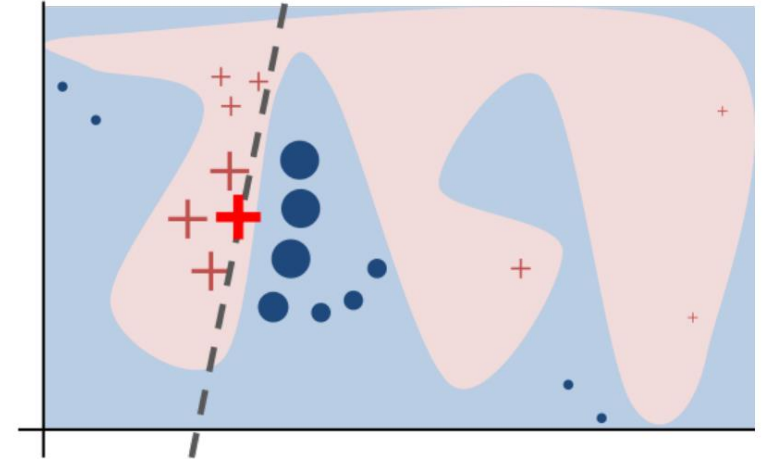
Minseok Yang (msyang0809@khu.ac.kr)  
Natural Language Processing Lab  
Department of AI, Kyung Hee University



## Part 6. Appendix

- **LIME (Marco Tulio Ribeiro et al., 2016)**

- Algorithm that can explain the predictions of any classifier or regressor in a faithful way, by approximating it locally with interpretable model



### Explanation produced by LIME

$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

### Formulation

- Explanation family
- Fidelity function
- Complexity measure

$G$   
 $\mathcal{L}$   
 $\Omega$

Minimize  $\mathcal{L}(f, g, \pi_x)$  while having  $\Omega(g)$  be low enough to be interpretable by humans

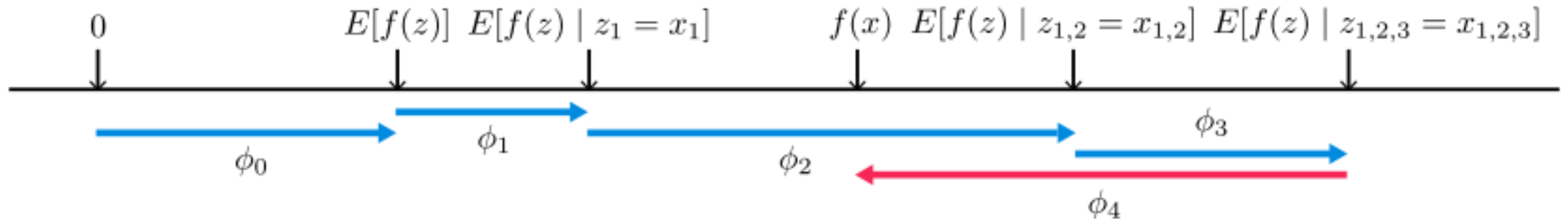
### Toy example to present intuition

- Blue/pink background
  - Black-box model's complex decision function (unknown to LIME)
- Red cross
  - Instance being explained
- Dashed line
  - The learned explanation that is locally (but not globally) faithful

# Appendix

## • KernelSHAP (Scott M Lundberg et al., 2017)

- Game theory results guaranteeing a unique solution apply to the class of additive feature attribution methods for model explanation

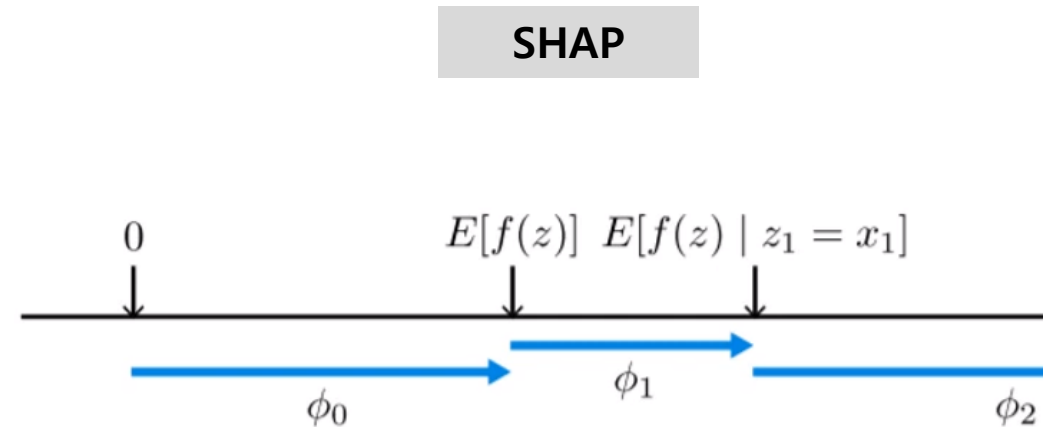
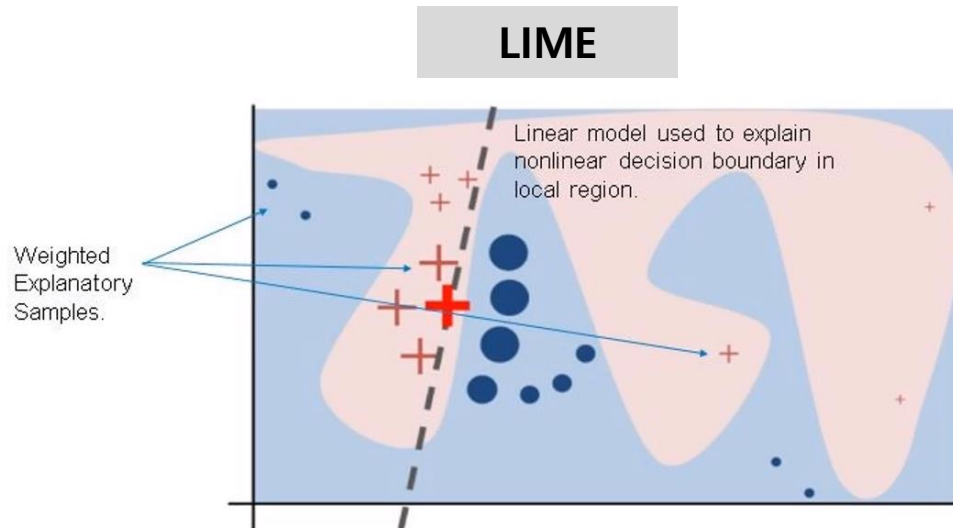


- SHAP (SHapley Additive exPlanation) values
  - Attribute to feature change in expected model prediction when conditioning on that feature
- To get from the base value  $E[f(z)]$ 
  - Predict if we did not know any features to the current output
- When the model is non-linear or the input features are not independent  $f(x)$ 
  - The order in which features are added to the expectation matters
  - SHAP values arise from averaging the  $\phi_i$  values across all possible orderings

## Part 6. Appendix

### • LIME & KernelSHAP

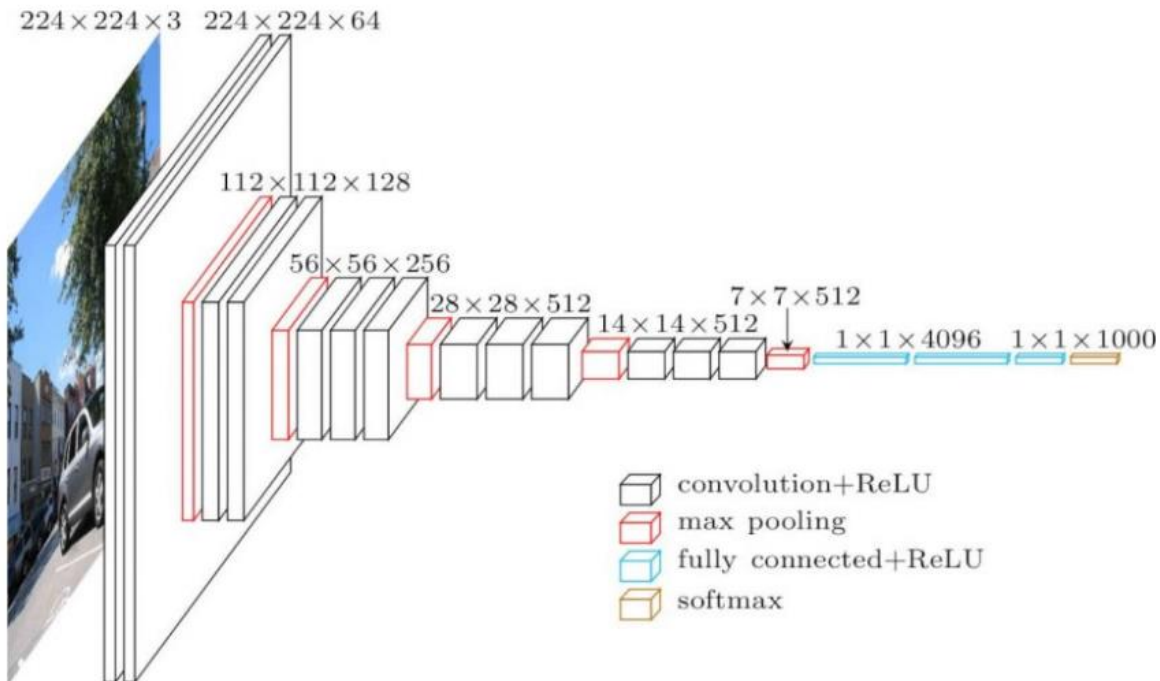
- LIME (Local Interpretable Model-agnostic Explanations)
  - Available on Tabular, Text, Image and even embedding Data
  - Instability because of variation on model explanation
- SHAP (Shapley Additive Explanations)
  - Based on Shapley Values as notion of game theory
  - Kernel SHAP ignore dependence among features



## Part 6. Appendix

- **VGG16 (Karen Simonyan et al., 2015)**

- Use very small  $3 \times 3$  receptive fields throughout the whole net, which are convolved with the input at every pixel
- All hidden layers are equipped with the rectification ReLU non-linearity



VGG16 Architecture

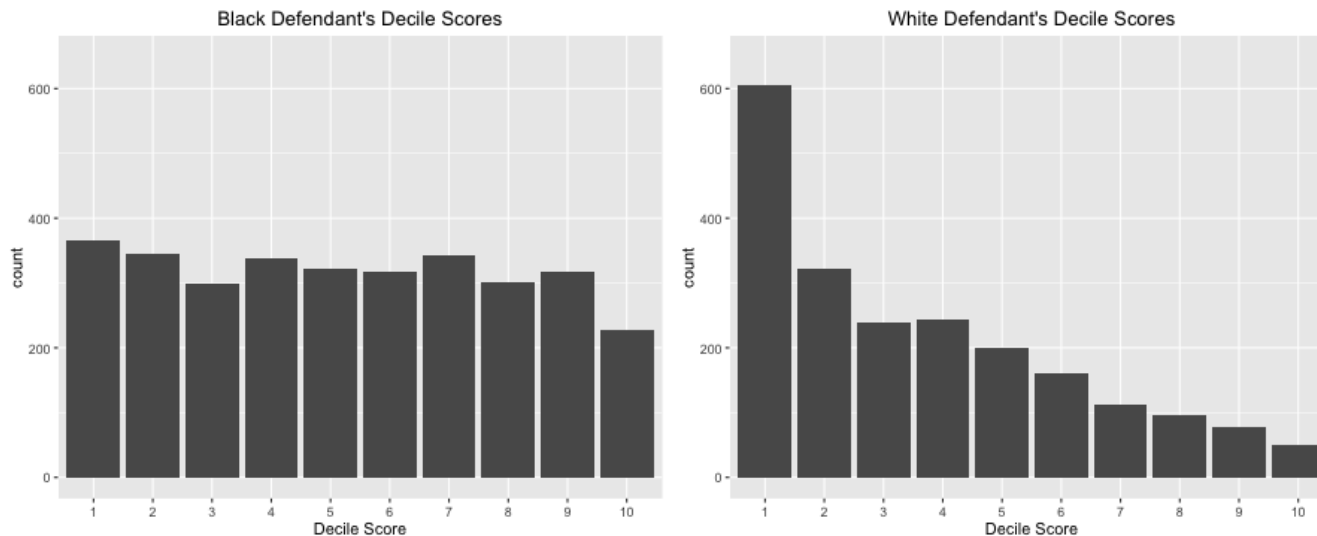
Table 3: ConvNet performance at a single test scale.

ConvNet config. (Table 1)	smallest image side		top-1 val. error (%)	top-5 val. error (%)
	train ( $S$ )	test ( $Q$ )		
A	256	256	29.6	10.4
A-LRN	256	256	29.7	10.5
B	256	256	28.7	9.9
C	256	256	28.1	9.4
	384	384	28.1	9.3
	[256;512]	384	27.3	8.8
D	256	256	27.0	8.8
	384	384	26.8	8.7
	[256;512]	384	25.6	8.1
E	256	256	27.3	9.0
	384	384	26.9	8.7
	[256;512]	384	<b>25.5</b>	<b>8.0</b>

## Part 6. Appendix

### • COMPAS

- Northpointe's tool, called COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)
  - Discover the underlying accuracy of their recidivism algorithm
  - Test whether the algorithm was biased against certain groups
- Data production
  - Looked at more than 10,000 criminal defendants in Broward County, Florida
  - Compared their predicted recidivism rates with the rate that actually occurred over a two-year period



#### Risk of Recidivism

- Black defendants were often predicted to be at a higher risk of recidivism than they actually were

#### Risk of Violent Recidivism

- To see this analysis result, please visit <https://github.com/propublica/compas-analysis>

# Appendix

## • German Credit

- Classifies people described by a set of attributes as good or bad credit risks
- Data form
  - Form provided by Prof. Hofmann, contains categorical/symbolic attributes
  - Form provided by Strathclyde University, is used for algorithms that need numerical attributes

### Cost matrix

	1	2
1	0	1
2	5	0

(1 = Good, 2 = Bad)

### Row

- The actual classification and the columns the predicted classification

### Data Sample

- It is worse to class a customer as good when they are bad (5), than it is to class a customer as bad when they are good (1)

# Appendix

- **MNIST (LeCun et al., 1998)**
  - Task: Image Classification
  - Handwritten number image data
    - Image & label from 0 to 9
    - Image: 28 X 28 pixel
    - Training Data: 60,000 / Test Data: 10,000





## Part 6. Appendix

- **ImageNet (J. Deng et al., 2009)**

- Task: Image Classification
- Object recognition, image classification and object localization
- 12 “subtrees”: mammal, bird, fish, reptile, amphibian, vehicle, furniture, musical instrument, geological formation, tool, flower, fruit

