

**ARTICLE TYPE**

# Visual Diagnostics of a Model Explainer – Tools for the Assessment of LIME Explanations

Katherine Goode\*<sup>1</sup> | Heike Hofmann<sup>1,2</sup><sup>1</sup>Department of Statistics, Iowa State University, Iowa, United States<sup>2</sup>Center for Statistics and Applications in Forensic Evidence (CSAFE), Iowa State University, Iowa, United States**Correspondence**

\*Corresponding author. Email: kgoode@iastate.edu

**Present Address**This is sample for present address text  
this is sample for present address text**Summary**

This is sample abstract text.

**KEYWORDS:**

LIME, black box models, interpretability, diagnostics

## 1 | INTRODUCTION

Machine learning algorithms have been proven to provide accurate predictions. However, this accuracy comes at a cost. With complicated algorithms, the ability to understand the mechanism in the data and the functional form of the model are lost. This has led for a desire to be able to interpret machine learning algorithms. One approach to remedy this problem would be to develop a model explainer. The objective of such an explainer is to understand and explain model performance.

Locally interpretable model explanations (LIME) is one such model explainer that was developed for this purpose (reference). While some model explainers are focused on understanding a model at the global level, LIME was designed to explain the model for a single prediction of interest. Conceptually, LIME fits a simple interpretable model (the explainer) that is meant to capture the behavior of the complex "black box model" in a local region around the prediction.

Predictive models are used in both regression and classification settings. For this paper, we will focus on the classification setting, and in particular, we will only discuss the case with a dichotomous response variable. In any case, when a predictive model is being developed, the data is usually divided into training and testing portions. The training data is used to fit the complex model, and the testing data is used to assess the model. The complex model can applied to the features in the testing data, and the resulting predictions can be compared to ground truth. Of interest during this assessment are the cases when the model is wrong. In the dichotomous response classification case, there are two ways in which the model can be wrong. The model can make a type I error in which..., or the model can make a type II error in which...

It does this by using the features from the training data to simulate a new dataset on which the simple model is fit. The complex model is applied to the simulated dataset to obtain predictions. The observations associated with predictions are used as the response variable in a ridge regression model with the the simulated features as the predictor variable with the highest weight given to observations closest to the prediction of interest. Feature selection is performed to identify the most important variables in the local region. A final ridge regression model is fit with the selected features, and the coefficients of the model are used to interpret the behavior of the complex model.

- explanation is also a prediction - how reliable is that explanation? - we will present some tools that could help answer that question - we will show an example using this tools (bullet data) - we will mention some possible insights about LIME that these tools have shown us 1. explainer model generally has very low  $R^2$  (probably due to binning) 2. "local" explanations are not local but are driven by the ("global") marginal distributions of covariates - Describe LIME including details on binning and linear regression in binned features and motivation for the binning

## 2 | METHODS

### 2.1 | LIME Algorithm

The LIME algorithm is presented in the original paper in a general context that includes cases for text classification and feature recognition. For this paper, we are only considering a situation with a binary categorical response variable and continuous feature variables. The procedure below is defined in this context. Furthermore, the implementation of LIME used in this paper is via the lime R package, and the procedure describes the methods used by the R package. This deviates a bit from the original procedure described by Riberiro, Singh, and Guestrin (2016). We attempted to highlight these deviations.

Let  $\mathbf{X}$  be an  $n$  by  $p$  data matrix with  $p$  features and  $n$  observations.

- $x = (x_1 \ x_2 \ \dots \ x_p) \in \mathbb{R}^p$ : original representation of an instance being explained (e.g. observations from a set of  $p$  continuous features for a specific case in the data)
- $x' \in \mathbb{R}^{p'}$ : vector for the interpretable representation of the instance being explained (e.g. a vector of indicator variables associated with the features chosen through feature selection indicating whether or not the observed value is in the bin created by LIME for the feature)
- $G$ : class of potentially interpretable models (e.g. linear models, decision trees, rule lists)
- $g$ : explanation model where  $g : \{0, 1\}^{p'} \rightarrow \mathbb{R}$  and  $g \in G$
- $\Omega(g)$ : measure of complexity of  $g$  (e.g. depth of a tree, number of non-zero coefficients in a linear model fit using LASSO)
- $f$ : model that is being explained where  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  (In classification  $f(x)$  is the probability that  $x$  belongs to a certain class)
- $\pi_x(z)$ : proximity measure between an instance  $z$  to  $x$  which defines a locality around  $x$
- $\mathcal{L}(f, g, \Pi_x)$ : the fidelity functions which is a measure of how unfaithful  $g$  is in approximating  $f$  in the locality defined by  $\pi_x$
- $\xi(x)$ : explanation produced by LIME where

$$\xi(x) = \arg \min_{g \in G} \mathcal{L}(f, g, \Pi_x) + \Omega(g)$$

(i.e. want to minimize  $\mathcal{L}(f, g, \Pi_x)$  and keep  $\Omega(g)$  low enough to be interpretable by humans)

### 2.2 | Application of LIME to Bullet Matching Data

#### 2.2.1 | lime R Package Perturbation Creation Methods

The LIME R package allows for the following four methods to sample the perturbations based on the distributions of the features from the training data.

- Equally Spaced Bins
- Quantile Bins

- Normal Approximation
- Kernel Density Approximation

The methods of equally spaced bins and quantile bins also allow the user to specify the number of bins. As of now, there are no recommendations or procedures provided for how to determine which method to use. By default, LIME uses four quantile bins. It was of interest to see how the explanations from LIME varied across the four sampling methods when applied to the bullet matching data. The LIME algorithm was applied to each prediction from the test data obtained from the 'rtrees' random forest model for each of the four sampling methods. Within the bin based sampling methods, the algorithm was applied for 2 to 6 bins. It was decided to only go up to 6 bins since the more bins used the more complex the explanation becomes.

## 2.2.2 | Proposed Bin Creation Methods

### 3 | DATA

### 4 | RESULTS

#### 4.1 | LIME Package Explanations Dependent on Sampling Method

In order to assess the LIME explanations created using different sampling methods, it was of interest to compare the top three features chosen as the important predictors by lime within a case from the test data across the different sampling methods. Figure ... is a heatmap showing the top feature chosen by lime for each of the cases in the test data and different bin based sampling methods. The rows represent the cases in the test data, and the columns represent the sampling methods. There are twenty methods included in the plot. These include the equally spaced bins and quantile bins from the lime package and the random forest score tree based bins and the same source tree based bins proposed in this paper. The rows are faceted by the test set and whether or not the observation is a match or not. The columns are faceted by these methods, and the columns within a facet represent the number of bins. Each method has 2 to 6 bins. The colors represent the top feature chosen by lime.

The variables of ccf and cms immediately show up as common variables chosen across all of the sampling methods. However, the patterns across the number of bins within the sampling methods are different. When equally spaced bins are used, the top feature chosen is consistent across all cases within a number of bins category. For example, ccf is almost always chosen (change to actual number) with 2 equally spaced bins, matches is almost always chosen with 3 equally spaced bins, and non\_cms is always chosen for the nonmatches with 5 and 6 equally spaced bins. This shows that with the bullet matching data, the top feature chosen with equally spaced bins is an artifact of the number of bins used. With equally spaced bins, this figure suggests that LIME is providing global explanations as opposed to local explanations. It would be preferable that the top feature chosen was more consistent across the number of bins and more variable across the cases. This would suggest that the top feature chosen is dependent on the feature values associated with a particular case and not just on which feature is the best explainer when  $b$  number of bins are used.

## ACKNOWLEDGMENTS

This is acknowledgment text [1]. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here. This is acknowledgment text. Provide text here.

## Author contributions

This is an author contribution text. This is an author contribution text. This is an author contribution text. This is an author contribution text. This is an author contribution text. This is an author contribution text.

