

## 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

To check on when editing writing:

- Generally: avoid 'can be'. Replace by 'is'.
- References like 'they', 'them' ... replace them by repeating the noun you are referring to to avoid any kind of ambiguity
- words to go back and make sure I am not using too often: ability, produce, understand

## 1 | INTRODUCTION

In the field of statistics, there are two main uses for models: inference and prediction. Machine learning models are often used for the latter purpose. While these models have been proven to perform well in a wide range of prediction problems, their accuracy comes at the cost of interpretability. They are often complex algorithms that are good at identifying patterns in the data but lack a functional form, so it not possible to directly interpret them. As a result, they produce predictions that lack an explanation, which has earned them the reputation of being “black box models”.

The ability to interpret a model serves multiple purposes. When a model is first fit, interpretations help to diagnose the model. Knowing which variables influence a prediction makes

it possible to determine if the predictions are based on reasonable variables. If not, appropriate adjustments to the model can be made. After a model has been diagnosed, the explanations for the predictions are used to understand the underlying mechanism that produced the data and provide an argument for why the predictions should be trusted. In some areas of application, this ability to be able to explain the results from a model is critical.

The forensics sciences and the health care industry are two areas that have benefited from machine learning but have an obvious need for explainable predictions. For example, Hare et al. [2] discuss the use of a random forest model in the forensics sciences to determine whether two bullets were fired from the same gun. Yu et al. [4] describe how healthcare has advanced through the use of machine learning in ways such as using neural networks to automate medical image diagnoses and implementing Bayesian networks to predict clinical outcomes. In both of these fields, the decisions made based on the results of the machine learning models can greatly impact people's lives. If it is not possible to explain the output from a model, it becomes questionable whether to rely on the predictions to make important decisions.

As an example, suppose that forensics examiners are trying to determine if a bullet used in a crime was fired from the gun that belongs to a suspect in the investigation. A random forest model could be used to produce a probability that the bullet was fired from the gun under question. Even if the model returns a high probability of a match, it will be difficult for a

jury to trust the results from the model when deciding whether to convict the suspect without the forensics examiners having the ability to explain which factors in the model led to the high probability.

The European Union has taken the desire to provide explanations from machine learning models a step further. In May 2018, the General Data Protection Regulation (GDPR) went into effect<sup>1</sup>. The regulation includes a policy that gives the right for individuals affected by decisions made via automated algorithms to understand the logic behind the decisions being made. As Goodman and Flaxman [1] point out, this regulation has a great effect on the way that machine learning algorithms are used to make decisions. The wording of the policy leaves room for interpretation, but the authors believe that, at a minimum, the regulations will require that an explanation of how the input variables relate to the predictions be provided.

As a result of the movement to produce explainable predictions, an area of research has emerged focusing on developing ways to explain output from machine learning algorithms. One approach that is being taken is to develop model explainers that provide insight into the performance of the complex model (reference). A model explainer is a method that is separate from the model but uses the model and output from the model to shed light on the process that the model goes through to produce predictions. LIME (local interpretable model-agnostic explanations) is one such model explainer [3].

While some model explainers are focused on understanding a model at the global level, LIME is designed to explain the model locally for a single prediction of interest. Additionally, LIME was designed to work with any model and to produce easily understandable explanations. XXX a citation might be good here for local and agnostic Conceptually, LIME fits a simple interpretable model, referred to as the explainer, that is meant to capture the behavior of the complex “black box model” in a local region around a prediction of interest. The interpretation of the simple model is then used to explain the variables that most influenced the prediction made by the complex model.

At a conceptual level, explainer models add another level of complexity to predictive models: in trying to explain the black box model, a simpler explainer model is added. To be able to trust in the explanation it is imperative to check that the explainer model is reliable and does not over-simplify the black-box model. It is important to keep in mind that the explanation produced by the explainer is not a direct interpretation of the complex model. Instead, it is based on a model that is trying to mimic a part of the complex model. Thus, it is necessary to ask the question, “How reliable is the explanation?”.

One way to answer this question is to understand how well the explainer model approximates the complex model.

The current implementations of LIME use a linear regression model as the explainer model (references to R and Python packages). These implementations are relying on the assumption that the relationship between the complex model predictions and the features is linear in a local region. It is important to assess this assumption in the local region of interest in order to trust the explanations produced by LIME. Additionally, the implementations offer various input settings, but little work has been done to provide advice on how to specify the settings in practice (look more into this). An assessment of the explanations could also help to determine which settings produce the explainer model with the best approximation of the complex model.

In this paper, we will present some visualization tools to assess the explanations from LIME. While predictive models are used in both regression and classification settings, we will focus on the classification setting for this paper. For additional simplicity, we will only discuss the case with a dichotomous response variable, but we believe that the work is adaptable for other situations.

Section 2 provides some background on explainer models that have been developed, a description of the LIME algorithm, and the current implementations of LIME. In Section 3, we discuss ways to assess the LIME explainer model and introduces our visualization tools. Section 4 demonstrates the use of our diagnostic tools with a logistic regression model fit to the iris data ( ? urgh. but we can decide on that later ) and a random forest model fit to a forensics bullet matching dataset. To conclude, Section 5 reviews the importance of assessing the LIME explainer model, discusses how the results from the example provide possible insights into the workings of LIME, and suggests future research directions.

## 2 | BACKGROUND

Just putting this here for now - needs lots of work or I may even remove it: Accepted best practice for all predictive models is that during the model fitting process, the data is 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. This is done to prevent overfitting and to get a more accurate assessment of the prediction error of the model. The complex model is then applied to the features in the testing data, and the resulting predictions are compared to ground truth. Of particular 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

<sup>1</sup><https://eugdpr.org/the-regulation/gdpr-faqs/>

wrong. The model can make a type I error in which..., or the model can make a type II error in which...

## 2.1 | Model Explainers

## 2.2 | Overview of LIME

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 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.

(a good place to include an image to help with this explanation)

## 2.3 | 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 Ribeiro, 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)

## 2.4 | Implementation of LIME in R

(section describing the procedure that is currently used by the LIME R package)

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.

## 3 | METHODS

Ways to understand if the LIME explainer is doing a good job:

- this figure suggest 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.

**How to cite this article:** Goode K., H. Hofmann, (2019), title, *journal*, volume, number and page.

