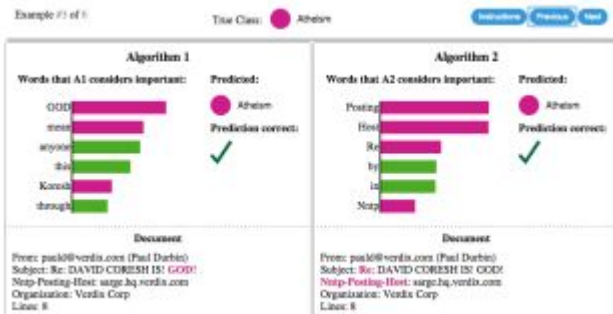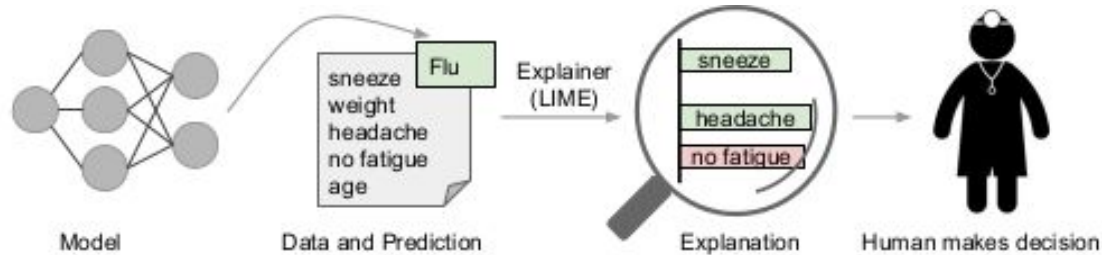# LIME

"Why Should I Trust You?"
Explaining the Predictions of Any Classifier
KDD '16

# Introduction

- With the increase in ML in tech, human should be able to trust the model.
- Trusting comprises of:
  - Trusting prediction
  - Trusting the model as a whole
- LIME: is an algorithm that can explain predictions of any classifier.
- SP-LIME: LIME with submodule optimization.

# Example



- With the help of which, humans with more domain knowledge of the problem will be able to collaborate.
- To ensure that features such as PARENTID do not contribute to the classification.

# Desired Characteristics for Explainers

- They must be interpretable
  - Explanations should be easy to understand
- Local fidelity
  - Prediction should be locally faithful
  - Model should behave similarly around the vicinity of the instance.
- Model-agnostic
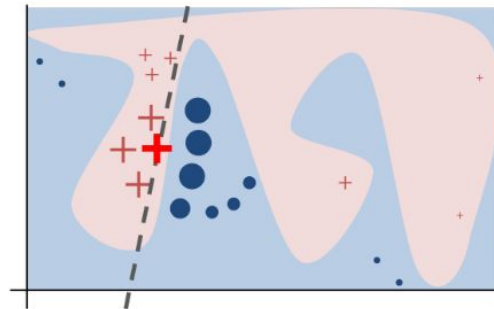  - Should be able to explain any model



**Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function $f$ (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using $f$, and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.**

# Basics - Interpretable Data Representation

- Explainer should use only interpretable representation
- Example
  - for text classification is a binary vector indicating the presence or absence of a word,
- This does not restrict the model, only the explainer.

# Loss function

- $\Omega(g)$ Is the measure of complexity of Explainer g
  - Ex, depth of decision tree.
  - Responsible for interpretability
- Model being explained is denoted by $f : \mathbb{R}^d \to \mathbb{R}$.
- $\pi_x(z)$ Is the proximity measure between z and x
  - Responsible for ensuring local fidelity
- $\mathcal{L}(f, g, \pi_x)$ Measures how unfaithful g is in approximating f in locality $\pi_x(z)$

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

Loss function

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) \left( f(z) - g(z') \right)^2$$

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

# Algorithm

**Algorithm 1** Sparse Linear Explanations using LIME

**Require:** Classifier $f$, Number of samples $N$
**Require:** Instance $x$, and its interpretable version $x'$
**Require:** Similarity kernel $\pi_x$, Length of explanation $K$

$\quad \mathcal{Z} \leftarrow \{\}$
$\quad$ **for** $i \in \{1, 2, 3, ..., N\}$ **do**
$\quad\quad z_i' \leftarrow sample\_around(x')$
$\quad\quad \mathcal{Z} \leftarrow \mathcal{Z} \cup \langle z_i', f(z_i), \pi_x(z_i) \rangle$
$\quad$ **end for**
$\quad w \leftarrow$ K-Lasso$(\mathcal{Z}, K)$ $\quad \triangleright$ with $z_i'$ as features, $f(z)$ as target
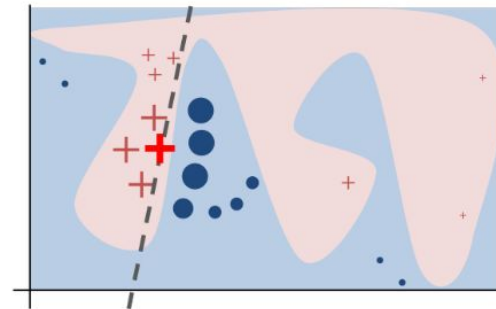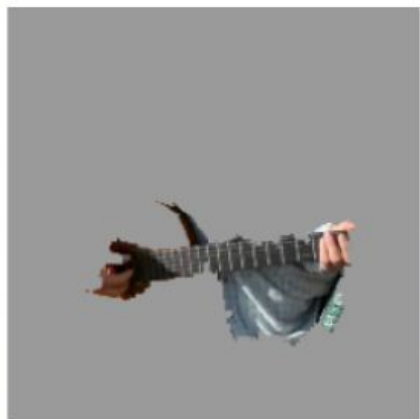$\quad$ **return** $w$



Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function $f$ (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using $f$, and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.
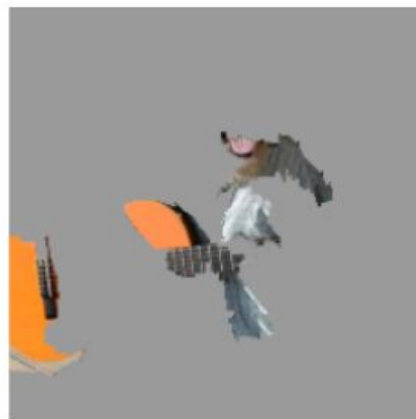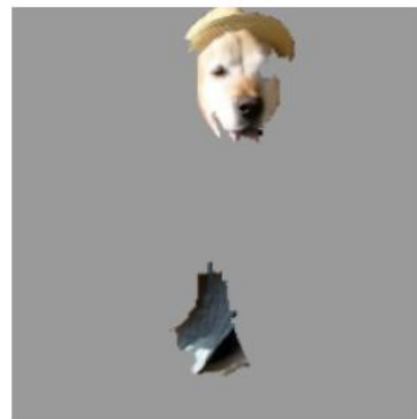
# Example - Image



(a) Original Image    (b) Explaining *Electric guitar*    (c) Explaining *Acoustic guitar*    (d) Explaining *Labrador*

Figure 4: Explaining an image classification prediction made by Google's Inception neural network. The top 3 classes predicted are "Electric Guitar" ($p = 0.32$), "Acoustic guitar" ($p = 0.24$) and "Labrador" ($p = 0.21$)

# SUBMODULAR PICK FOR EXPLAINING MODELS

**Algorithm 2** Submodular pick (SP) algorithm

**Require:** Instances $X$, Budget $B$
   **for all** $x_i \in X$ **do**
      $\mathcal{W}_i \leftarrow \textbf{explain}(x_i, x_i')$          ▷ Using Algorithm 1
   **end for**
   **for** $j \in \{1 \dots d'\}$ **do**
      $I_j \leftarrow \sqrt{\sum_{i=1}^n |\mathcal{W}_{ij}|}$    ▷ Compute feature importances
   **end for**
   $V \leftarrow \{\}$
   **while** $|V| < B$ **do**        ▷ Greedy optimization of Eq (4)
      $V \leftarrow V \cup \text{argmax}_i \, c(V \cup \{i\}, \mathcal{W}, I)$
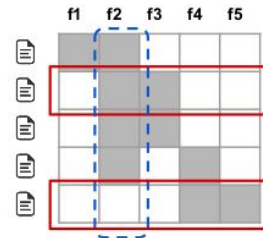   **end while**
   **return** $V$



Figure 5: Toy example $\mathcal{W}$. Rows represent instances (documents) and columns represent features (words). Feature f2 (dotted blue) has the highest importance. Rows 2 and 5 (in red) would be selected by the pick procedure, covering all but feature f1.

$$c(V, \mathcal{W}, I) = \sum_{j=1}^{d'} \mathbb{1}_{[\exists i \in V : \mathcal{W}_{ij} > 0]} I_j$$
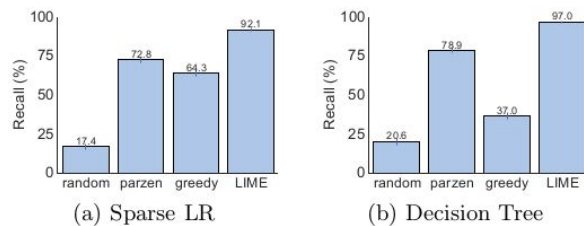
# Experiments



(a) Sparse LR     (b) Decision Tree

**Figure 6: Recall on truly important features for two interpretable classifiers on the books dataset.**
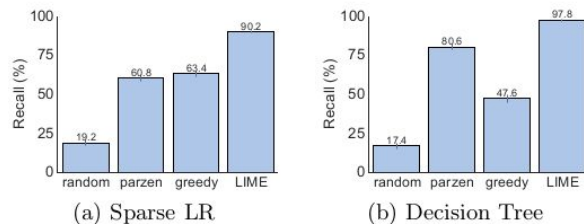


(a) Sparse LR     (b) Decision Tree

**Figure 7: Recall on truly important features for two interpretable classifiers on the DVDs dataset.**

# Experiment

- They consider a prediction trustworthy if it does not change with removal of untrustworthy feature.
- Untrustworthy features are selected randomly 25% of total features (identified by users)
- For Lime after removing untrustworthy features in explanation, if prediction changes then it is untrustworthy

Table 1: Average F1 of *trustworthiness* for different explainers on a collection of classifiers and datasets.

| | Books | | | | DVDs | | | |
|---|---|---|---|---|---|---|---|---|
| | LR | NN | RF | SVM | LR | NN | RF | SVM |
| Random | 14.6 | 14.8 | 14.7 | 14.7 | 14.2 | 14.3 | 14.5 | 14.4 |
| Parzen | 84.0 | 87.6 | 94.3 | 92.3 | 87.0 | 81.7 | 94.2 | 87.3 |
| Greedy | 53.7 | 47.4 | 45.0 | 53.3 | 52.4 | 58.1 | 46.6 | 55.1 |
| LIME | **96.6** | **94.5** | **96.2** | **96.7** | **96.6** | **91.8** | **96.1** | **95.6** |

# "Husky vs Wolf"- insight from LIME



(a) Husky classified as wolf    (b) Explanation

Figure 11: Raw data and explanation of a bad model's prediction in the "Husky vs Wolf" task.

|  | Before | After |
|---|---|---|
| Trusted the bad model | 10 out of 27 | 3 out of 27 |
| Snow as a potential feature | 12 out of 27 | 25 out of 27 |

Table 2: "Husky vs Wolf" experiment results.

# Conclusion

- Lime is introduced as an modular and extensible approach to faithfully explaining the predictions
- SP-Lime was also introduced which helped in selecting non redundant features.