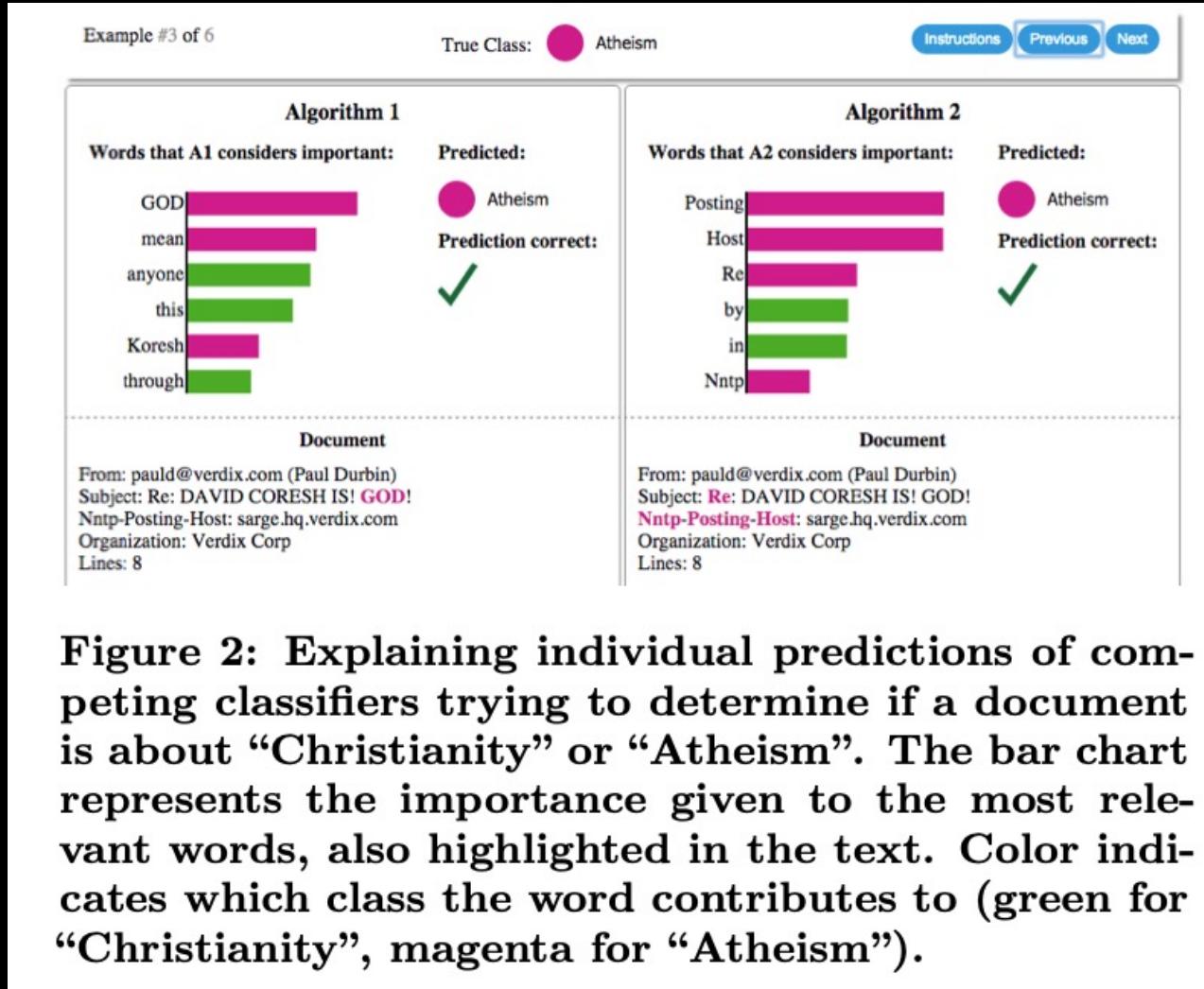


# Explainable AI: LIME & SHAP - I

Dive deep into the Math, intuition & code

AppliedAIcourse.com

Why does a model in a specific way  
for a given input ?



**Figure 2: Explaining individual predictions of competing classifiers trying to determine if a document is about “Christianity” or “Atheism”. The bar chart represents the importance given to the most relevant words, also highlighted in the text. Color indicates which class the word contributes to (green for “Christianity”, magenta for “Atheism”).**

most images from <https://arxiv.org/pdf/1602.04938.pdf>

Explainability  $\approx$  Feature importance  
and/or  
Feature contribution

Linear models: magnitude of weights

Tree based models: Feature importance via  
Entropy or Info Gain

Deep Learning models: Integrated Gradients

[https://www.tensorflow.org/tutorials/interpretability/integrated\\_gradients](https://www.tensorflow.org/tutorials/interpretability/integrated_gradients)

Can we have model-agnostic explainability

LIME (2016)

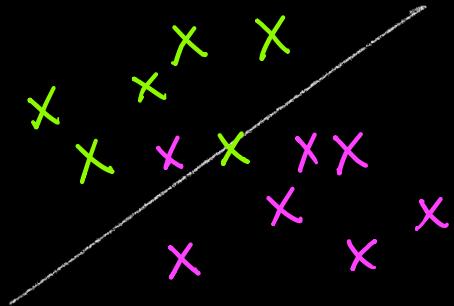
<https://arxiv.org/pdf/1602.04938.pdf>

Local interpretable model agnostic - explanations

Paper: "why should I trust you?"

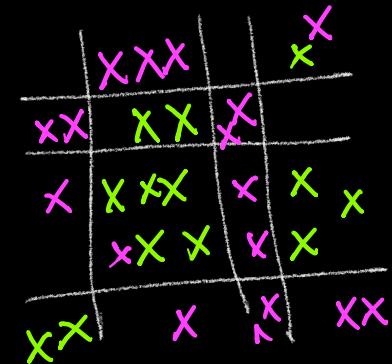
# Local vs Global Interpretability

Linear models:



→ Global

Tree based models:

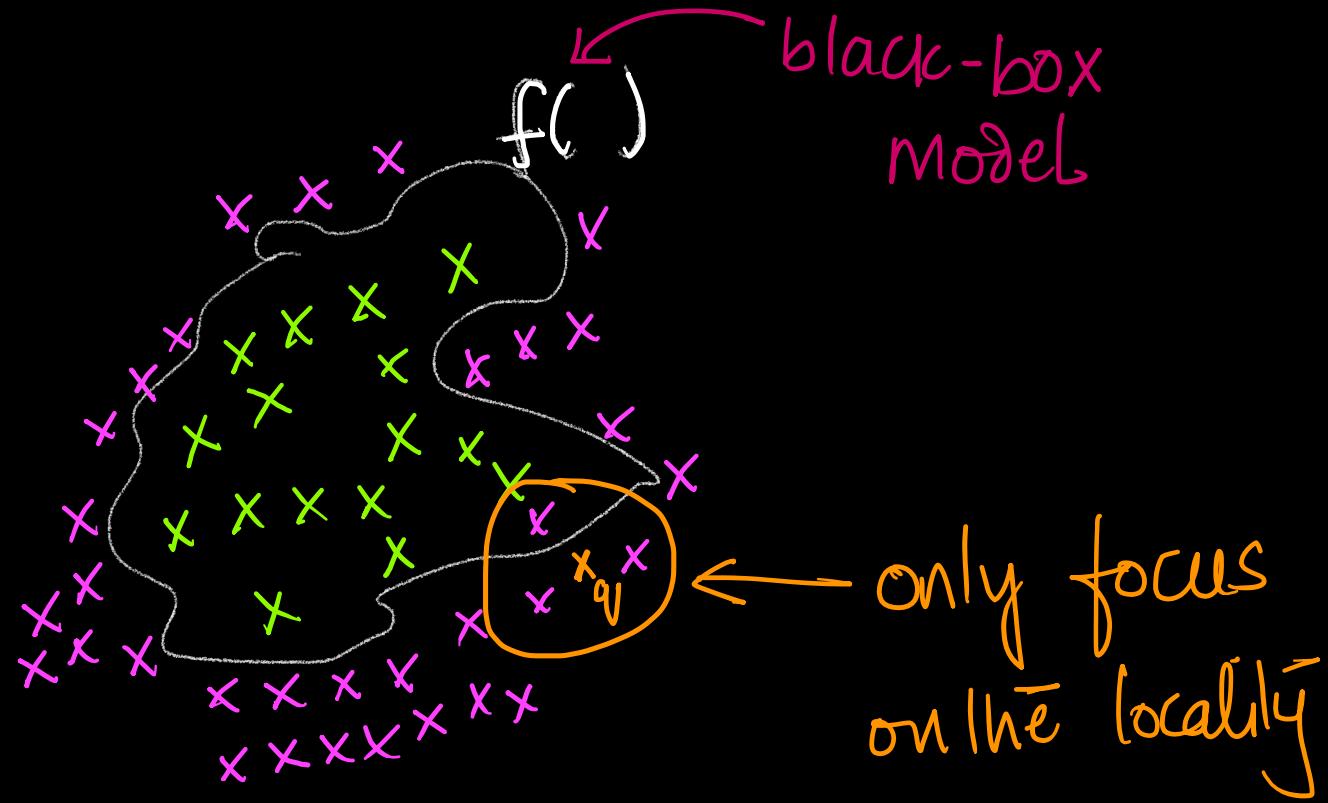


→ Global

Local-interpretability

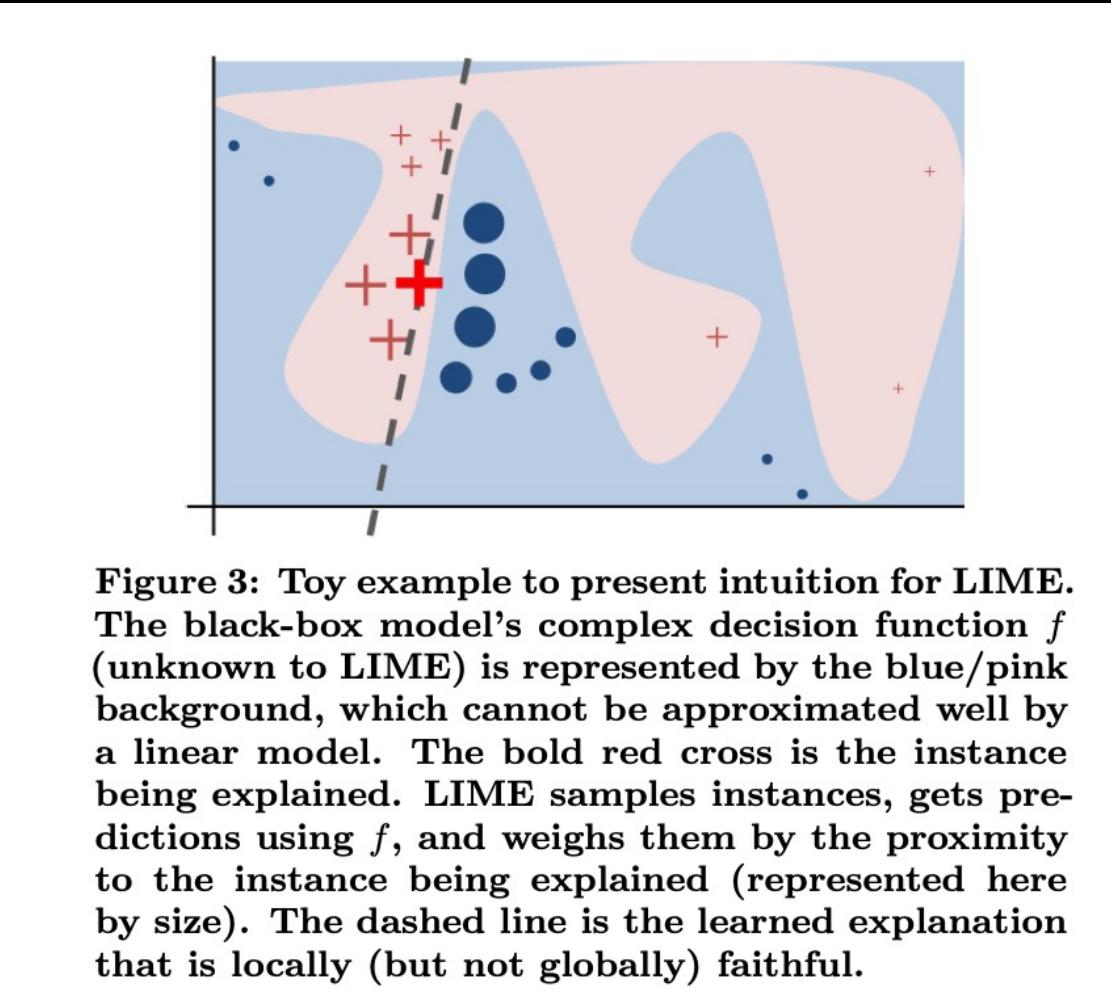
$f(x) = y \in \mathbb{R}$

- regression
- probability for classification



why is the model behaving in a specific way  
in the locality of  $z_q$ ?

Intuition:



Math & Algorithm:

$$x \in \mathbb{R}^d \rightarrow f(\cdot) \rightarrow y$$

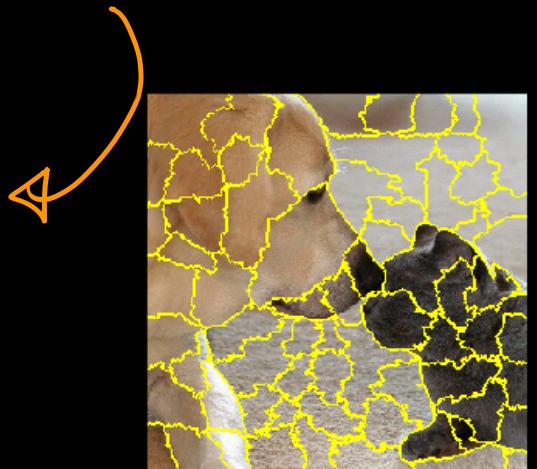
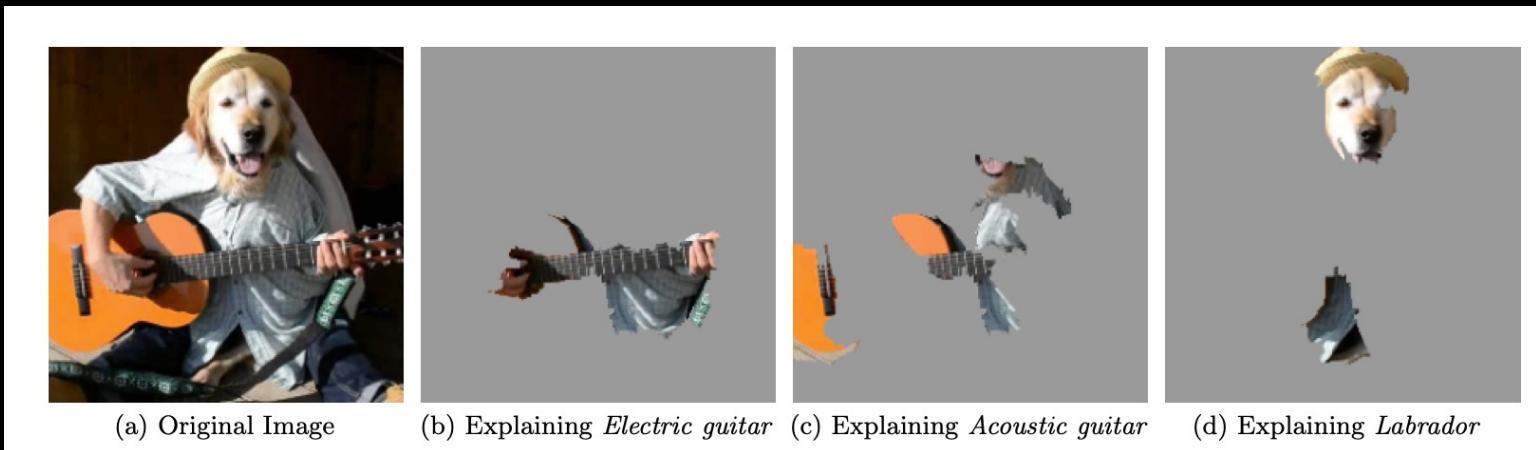
Black Box

$$x' \in \{0,1\}^{d'}$$

↳  $d'$ -dim interpretable representation as  
binary vector (presence/absence)

$$x' \in \{0, 1\}^{d'}$$

→ Hack: Limit to k-words  
presence | absence of a word (bow)  
presence | absence of an image patch  
or Super-pixel



Surrogate model :  $g$

model  $g \in G$  = class of easily interpretable models  
(L<sub>0</sub>-models, DTs)

$\Omega(g)$  = complexity of  $g$  (e.g. depth of DT / # of non-zero weights)

$g(x')$  and not  $g(x)$

↑ interpretable representation

↑  $\in \mathbb{R}^d$ , original feature vector

Proximity Function

$\pi_x(z)$  = proximity measure between  $x$  &  $z$

$\pi_x \rightarrow$  locality of  $x$

Local Fidelity:

$L(f, g, \pi_x)$ : How unfaithfully  $g()$  approximates  $f()$  in the locality defined by  $\pi_x$

$$\xi(x) = \min_{g \in G} L(f, g, \pi_x) + \mathcal{L}(g) \rightarrow \boxed{1}$$

$G$ : linear models

$\mathcal{L}$ : squared loss as  $f(x) \in \mathbb{R}$

$$\pi_x(z) = \exp \left\{ - \frac{\mathcal{D}(x, z)^2 / \sigma^2}{2} \right\} \rightarrow \text{exponential kernel}$$

↑  
dist

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

$$\xi(x) = \min_{g \in G} \mathcal{L}(f, g, \Pi_{\mathcal{X}}) + \boxed{\mathcal{L}(g)}$$

K-Lasso

K-Lasso = L<sub>1</sub>-regularized linear regression  
with only top K-features

---

**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  $K$

$\mathcal{Z} \leftarrow \{\}$

**for**  $i \in \{1, 2, 3, \dots, N\}$  **do**

$z'_i \leftarrow sample\_around(x')$

$\mathcal{Z} \leftarrow \mathcal{Z} \cup \langle z'_i, f(z_i), \pi_x(z_i) \rangle$

**end for**

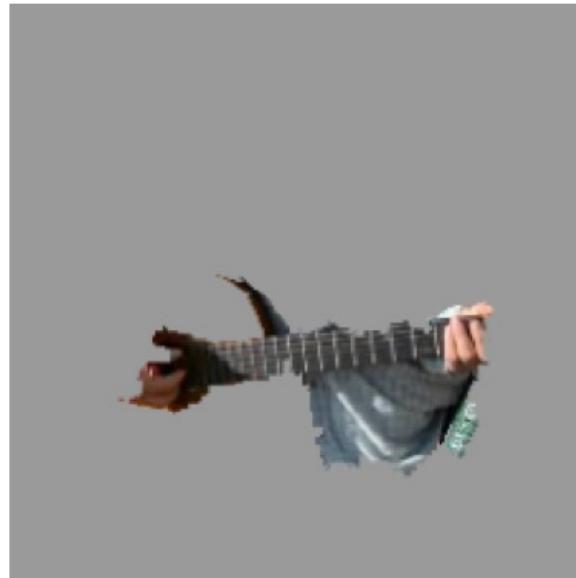
$w \leftarrow K\text{-Lasso}(\mathcal{Z}, K)$   $\triangleright$  with  $z'_i$  as features,  $f(z)$  as target

**return**  $w$

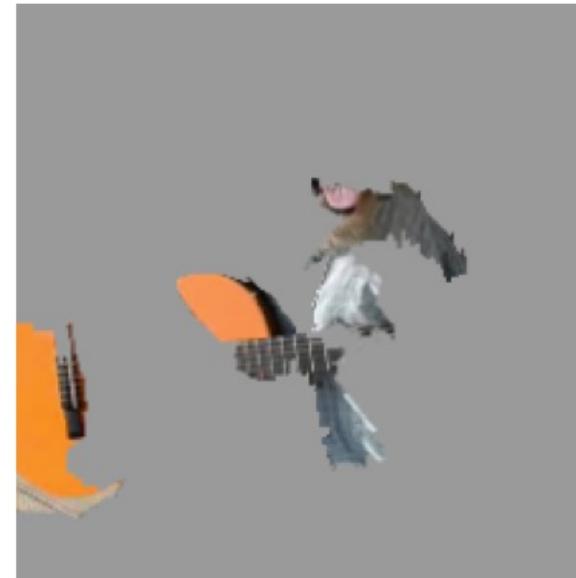
---



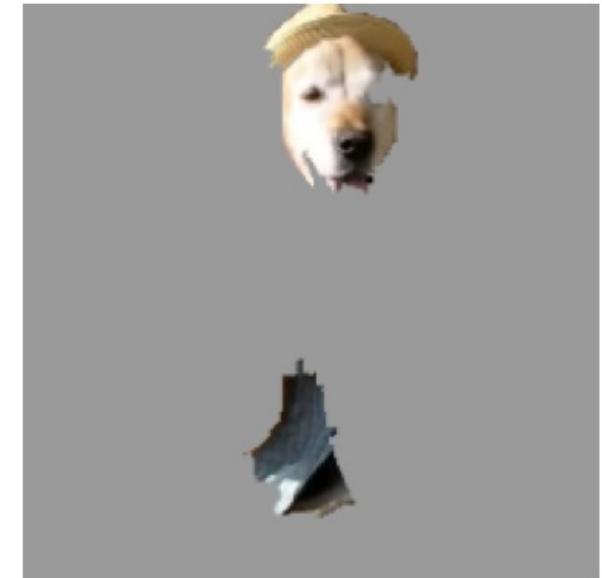
(a) Original Image



(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” ( $p = 0.21$ )**

## Advantages

- Works on Images, Text & Tabular data
- Lasso or low-depth trees for  $G \rightarrow$  highly interpretable
- $x$  need not be same as  $x'$ 
  - ↓  
may/may not be  
interpretable
  - ↓  
interpretable feature vector

## Drawbacks

- $\Pi_x(z)$ : What is neighborhood?
- Kernel width ( $\sigma$ ) → what is the right value  
(Try various values) ↗
- distance metric:  $\mathcal{D}(x, z)$  in high-dimensional data
- $x \rightarrow x' \in \{0, 1\}^{d'}$  interpretable features may not be straight forward  
e.g.: real-valued features → bucketing

## - Instability of explanations

<https://arxiv.org/pdf/1806.08049.pdf>

↳ simulated data: explanations of close points varied

↳ repeat the sampling could give different explanations

⇒ Lack of robustness

Alternative:

SHAP → uses Shapley values from Game Theory

→ concrete Mathematical guarantees

↓ Time-consuming

→ harder to interpret  $\Rightarrow$  more likely to misinterpret

Nice web-book ↴

<https://christophm.github.io/interpretable-ml-book/index.html>

[covers breadth well]

for depth → checkout the research papers

Code [very-simple] : <https://lime-ml.readthedocs.io/en/latest/>

Tabular ↴

<https://towardsdatascience.com/lime-how-to-interpret-machine-learning-models-with-python-94b0e7e4432e>

Image ↴

<https://towardsdatascience.com/interpreting-image-classification-model-with-lime-1e7064a2f2e5>

Text ↴

<https://towardsdatascience.com/what-makes-your-question-insincere-in-quora-26ee7658b010>

Next - Session : SHAP

Q&A