

# Explaining Black-Box Machine Learning Predictions

Sameer Singh

University of California, Irvine

# Machine Learning is Everywhere...



# Classification: Wolf or a Husky?



Machine  
Learning  
Model



**Wolf!**

# Classification: Wolf or a Husky?



Machine  
Learning  
Model



Husky!

# Classification: Wolf or a Husky?

Only 1 mistake!



Predicted: **wolf**  
True: **wolf**



Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**



Predicted: **wolf**  
True: **husky**

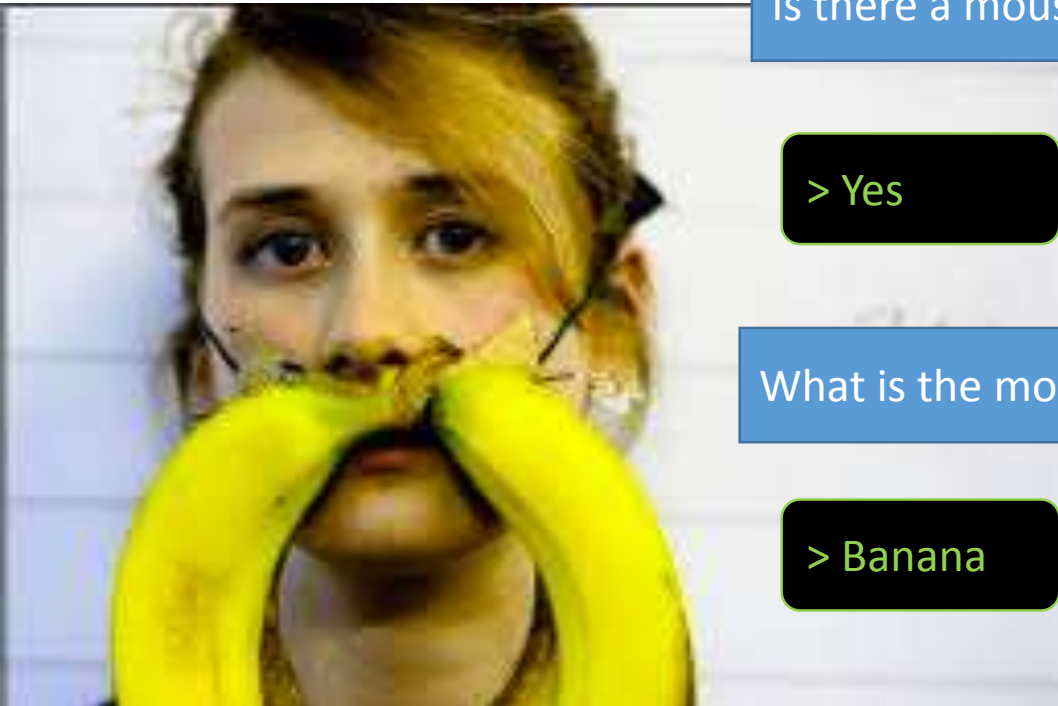


Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**

# More Complex: Question Answering

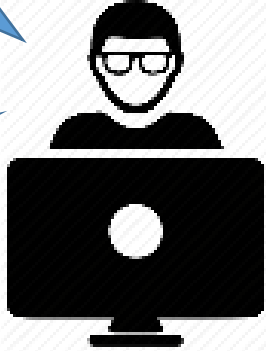


Is there a moustache in the picture?

> Yes

What is the moustache made of?

> Banana



The diagram illustrates a question-answering process. On the left, a photograph of a woman with a banana mustache is shown. To its right, a series of blue speech bubbles and green response boxes are connected by lines. The first speech bubble asks 'Is there a moustache in the picture?', followed by a green box with the answer '> Yes'. The second speech bubble asks 'What is the moustache made of?', followed by a green box with the answer '> Banana'. To the right of these exchanges is a black icon of a person wearing glasses and sitting at a computer, representing the system or user initiating the questions.

# Essentially black-boxes!

## Trust

How can we trust the predictions are correct?

How do we know they are not breaking regulations?

How do we avoid “stupid mistakes”?

## Predict

How can we understand and predict the behavior?

## Improve

How do we improve it to prevent potential mistakes?



# Classification: Wolf or a Husky?

Only 1 mistake!



Predicted: **wolf**  
True: **wolf**



Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**



Predicted: **wolf**  
True: **husky**



Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**

We've built a  
snow detector...



VIDEO

SLATE IN MOTION.

OCT. 14 2016 3:18 PM

# The Man Who Accidentally Adopted a Wolf Pup

It did not go well.

By *A.J. McCarthy*



10k



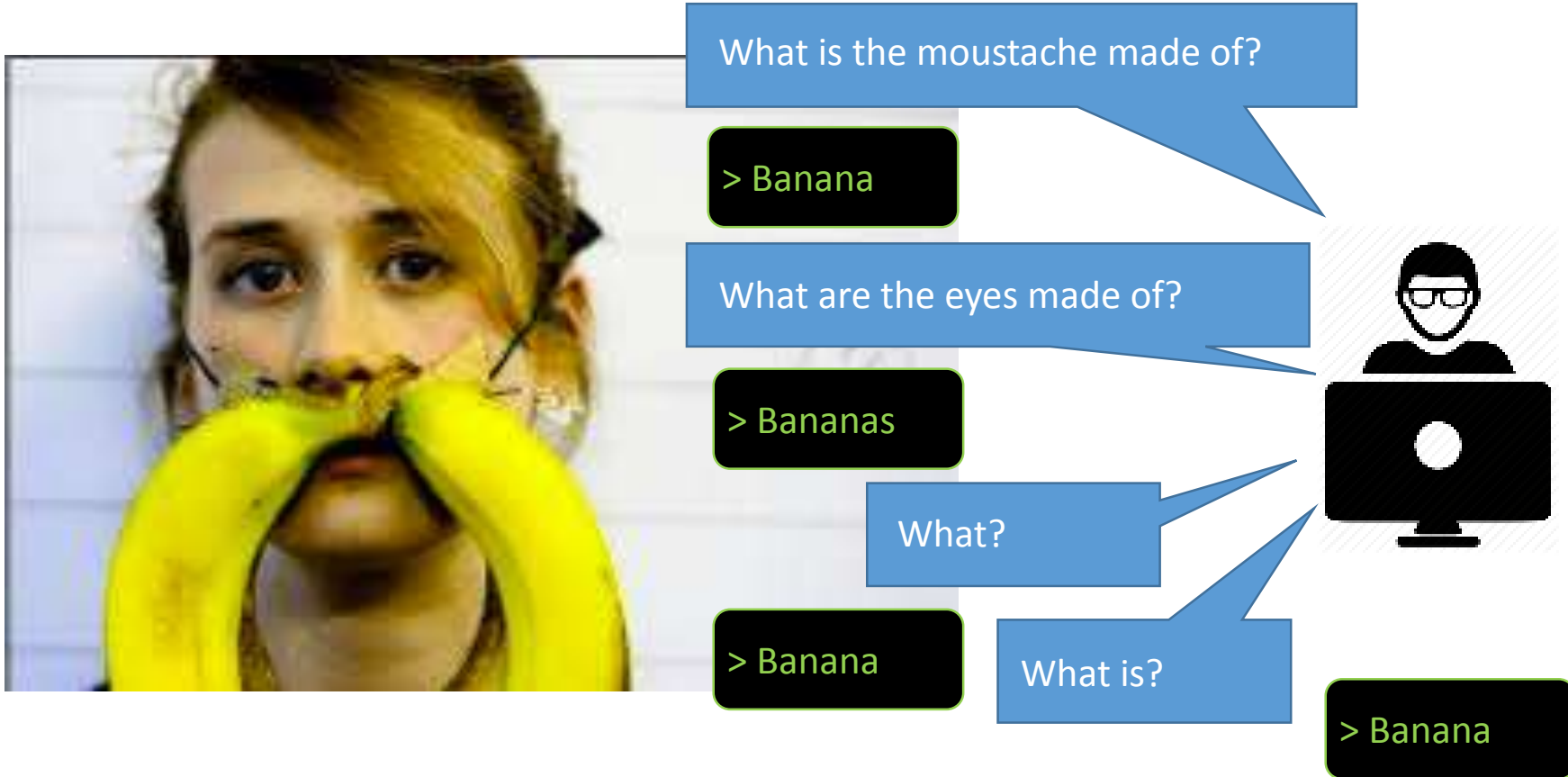
547



6



# Visual Question Answering



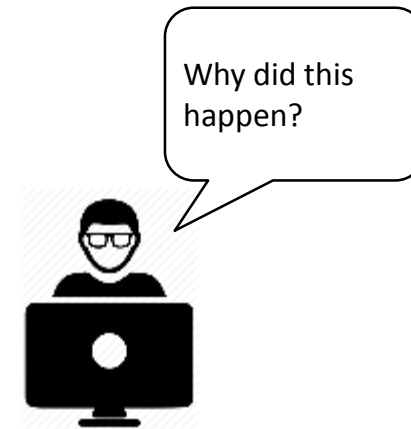
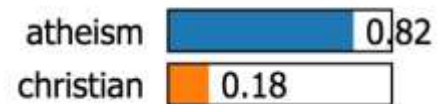
# Text Classification

From: Keith Richards  
Subject: Christianity is the answer  
NTTP-Posting-Host: x.x.com

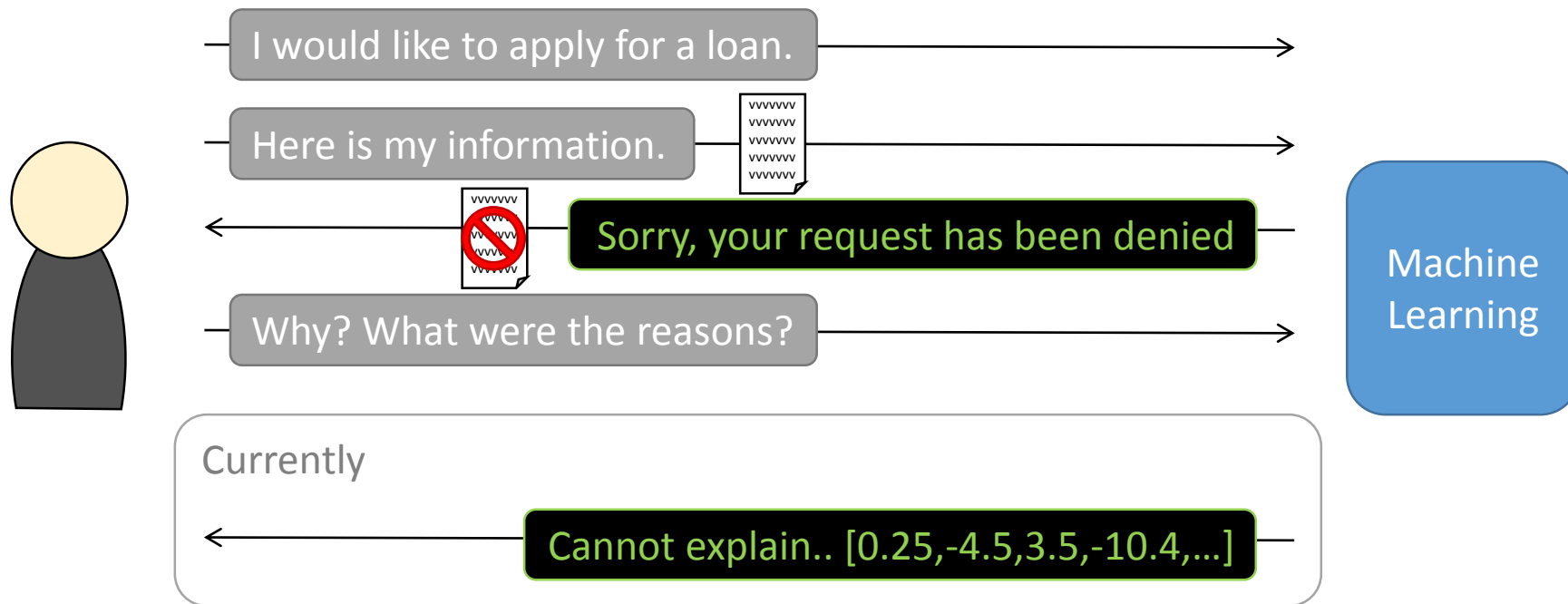
I think Christianity is the one true religion.  
If you'd like to know more, send me a note



Prediction probabilities



# Applying for a Loan

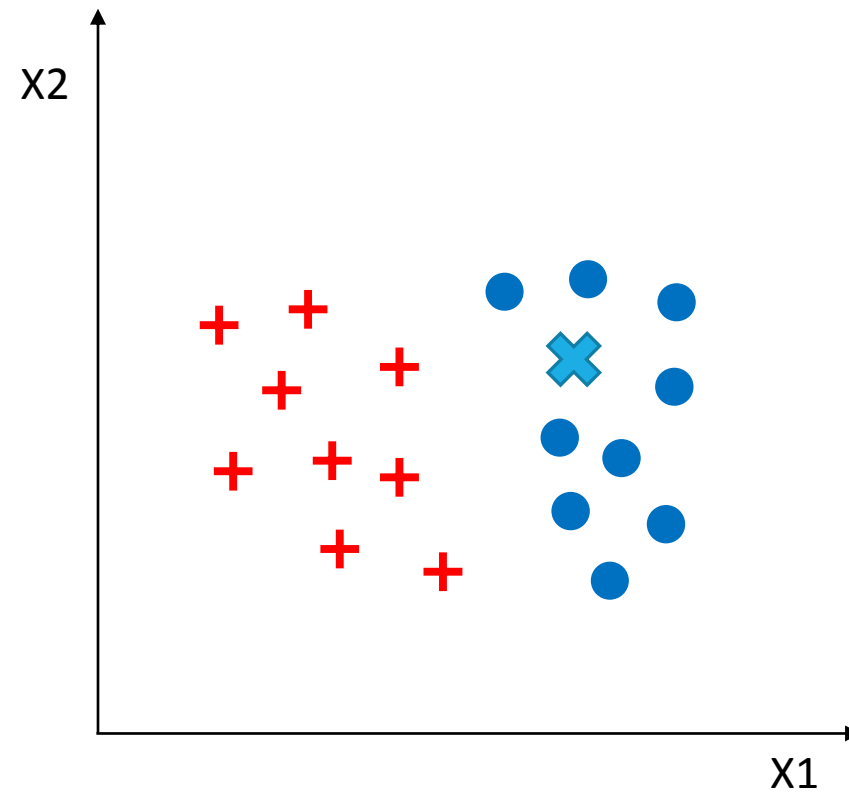


# How did we get here?

Big Data and Deep Learning

# Simple Data

---

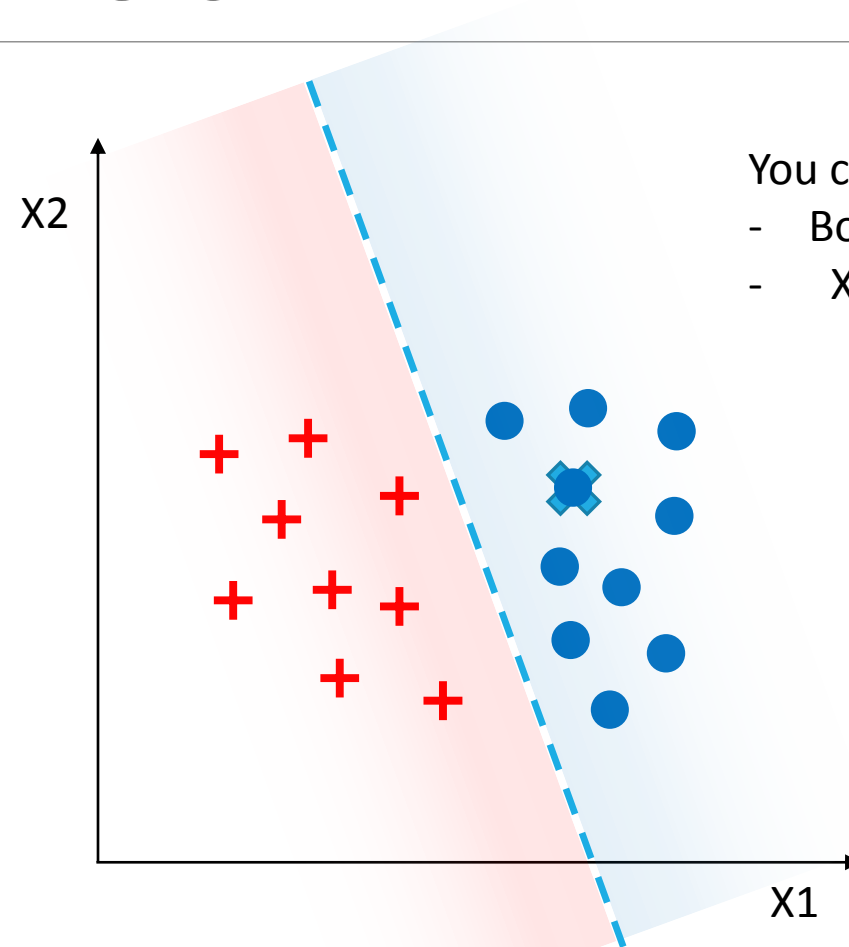


# Linear Classifiers

if:  $10X_1 + X_2 - 5 > 0$

otherwise

+

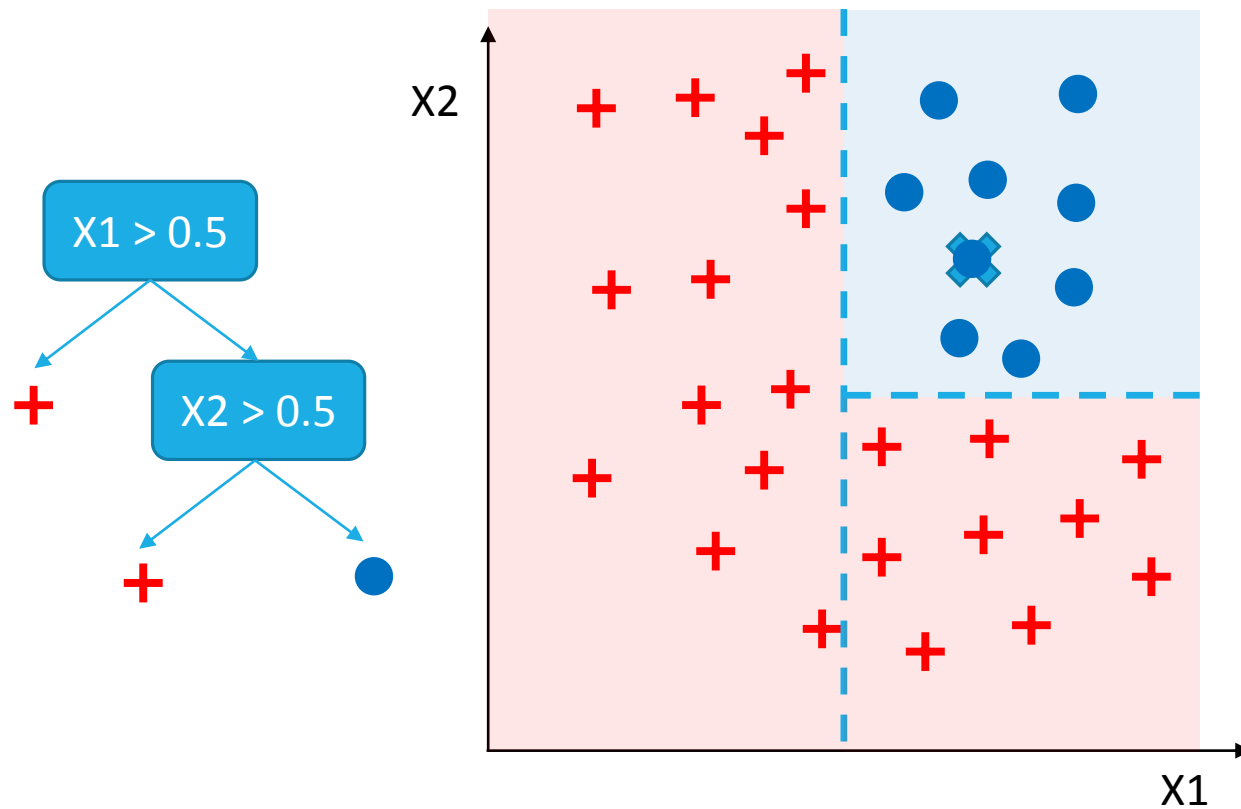


You can interpret it...

- Both have a positive effect
- $X_1 > X_2$



# Decision trees

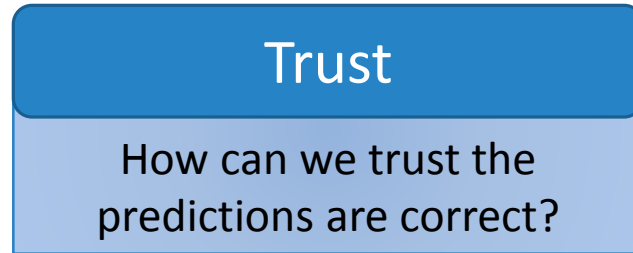


You can interpret it...

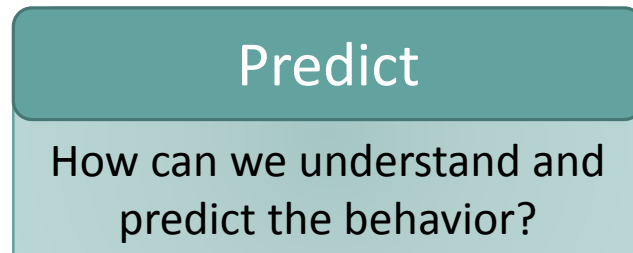
- $X_2$  is irrelevant if  $X_1 < 0.5$
- Otherwise  $X_2$  is enough

# Looking at the structure

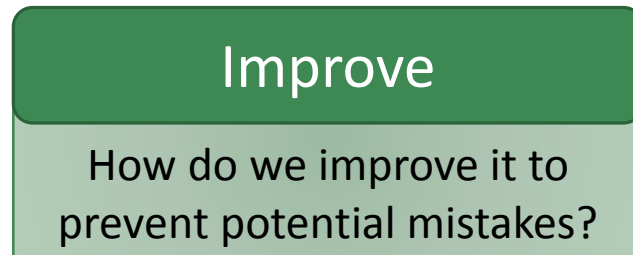
---



Test whether the structure agrees with our intuitions.



Structure tells us exactly what will happen on any data.



Structure tells you where the error is, thus how to fix it.

# Arrival of Big Data

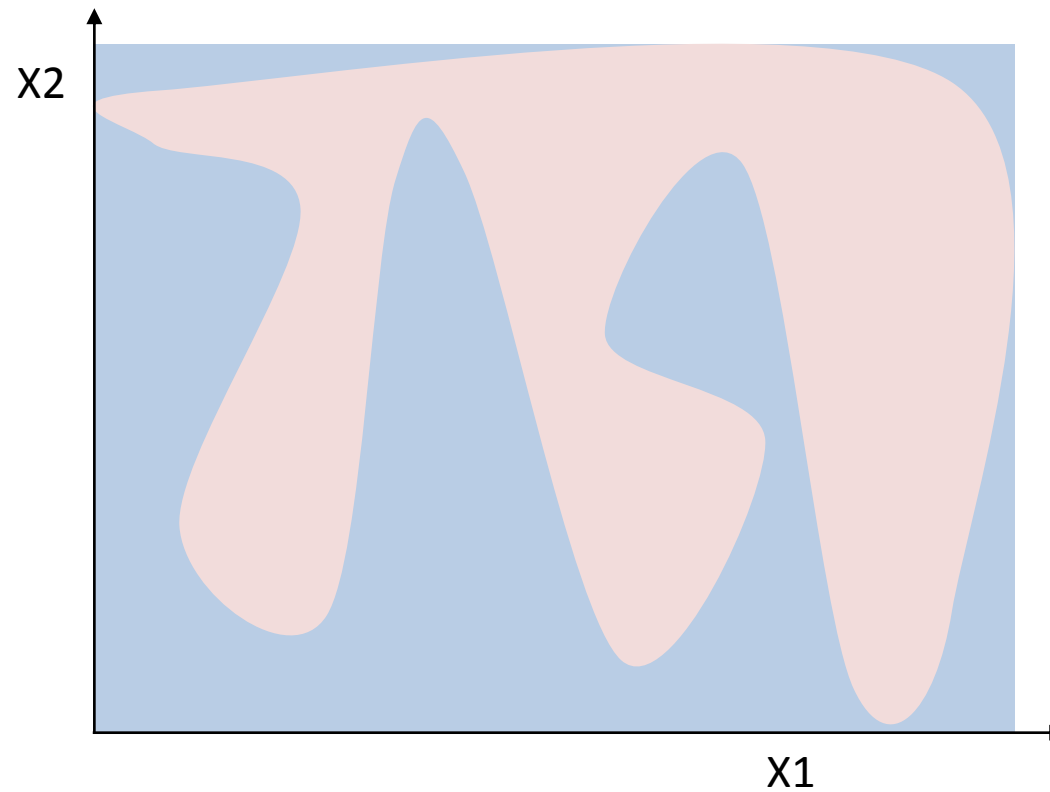
---

# Big Data: Applications of ML



# Big Data: More Complexity

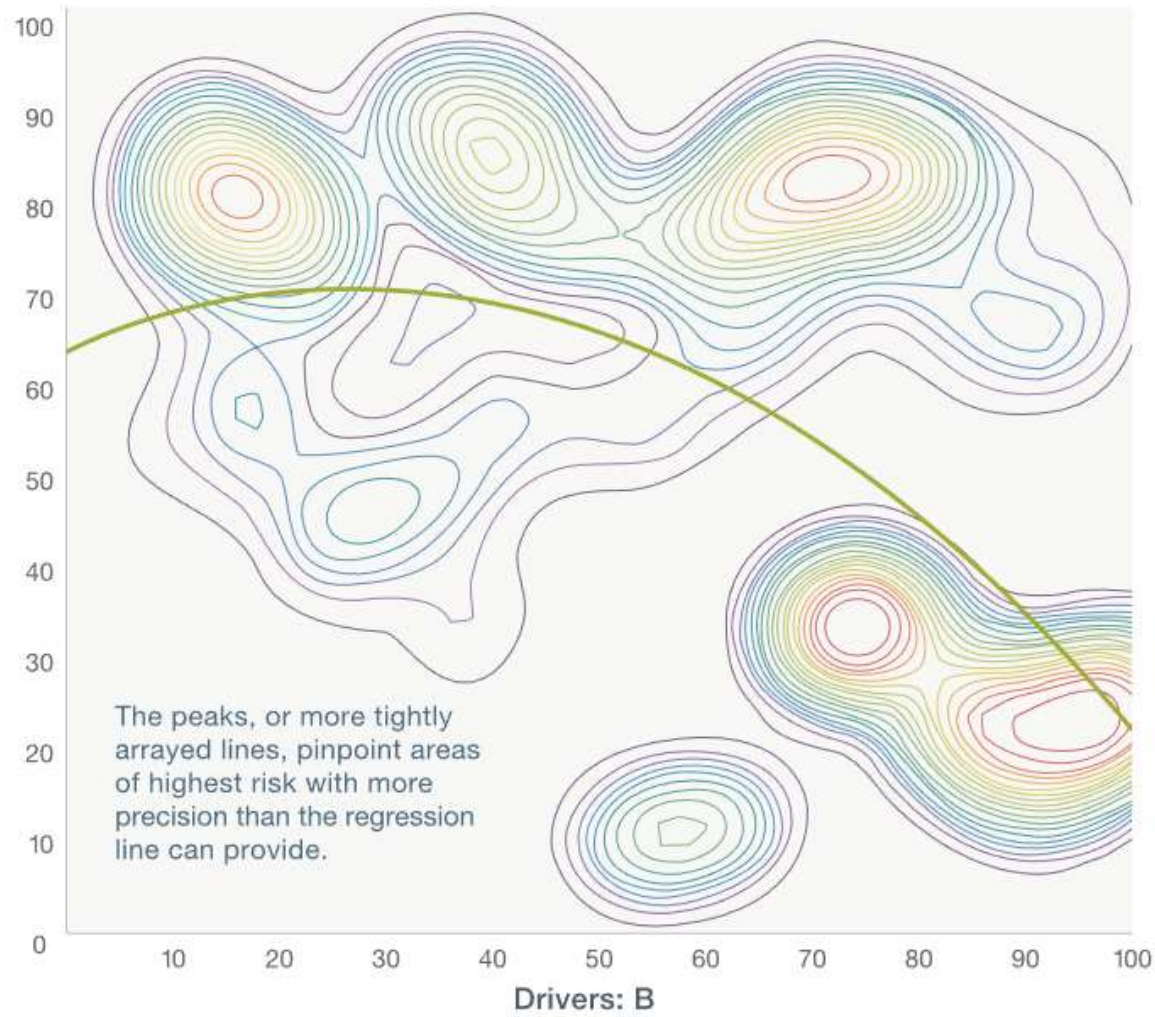
---



— Classic regression analysis

○ Isobar graph facilitated by machine learning: warmer colors indicate higher degrees of risk

Drivers: A



McKinsey&Company

# Big Data: More Dimensions

---

Savings

Income

Credit scores

Loan Amount

Past defaults

Recent defaults

Profession

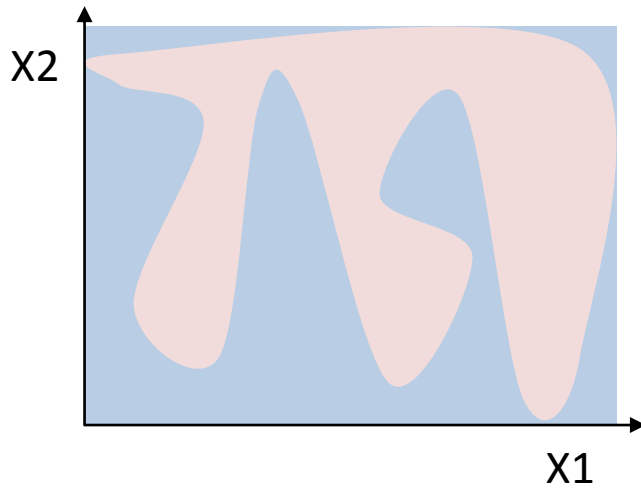
Age

Marital  
Status

This **easily** goes to hundreds

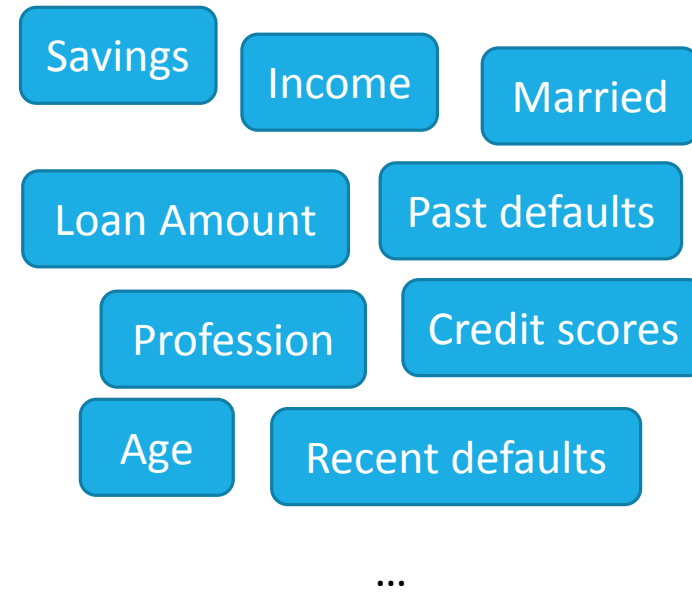
- Images: thousands
- Text: tens of thousands
- Video: millions
- ... and so on





Complex Surfaces

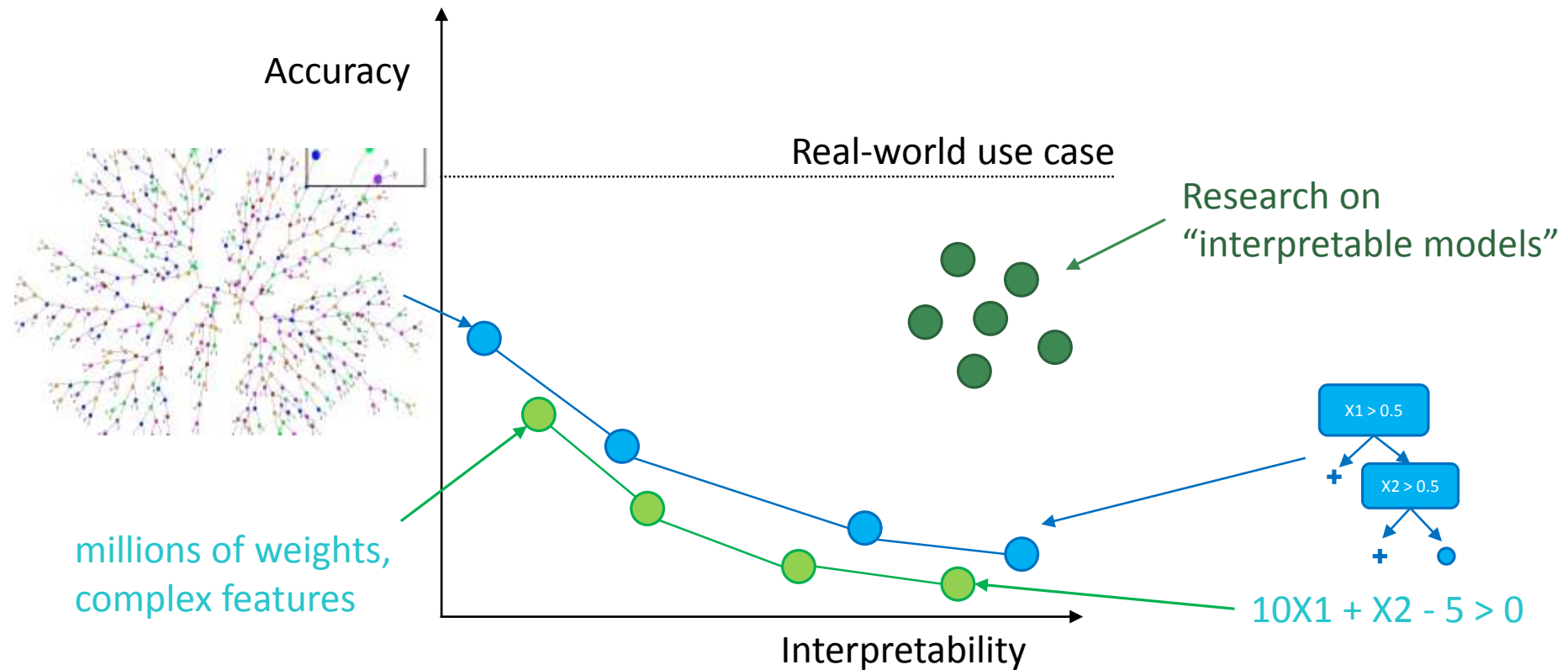
+



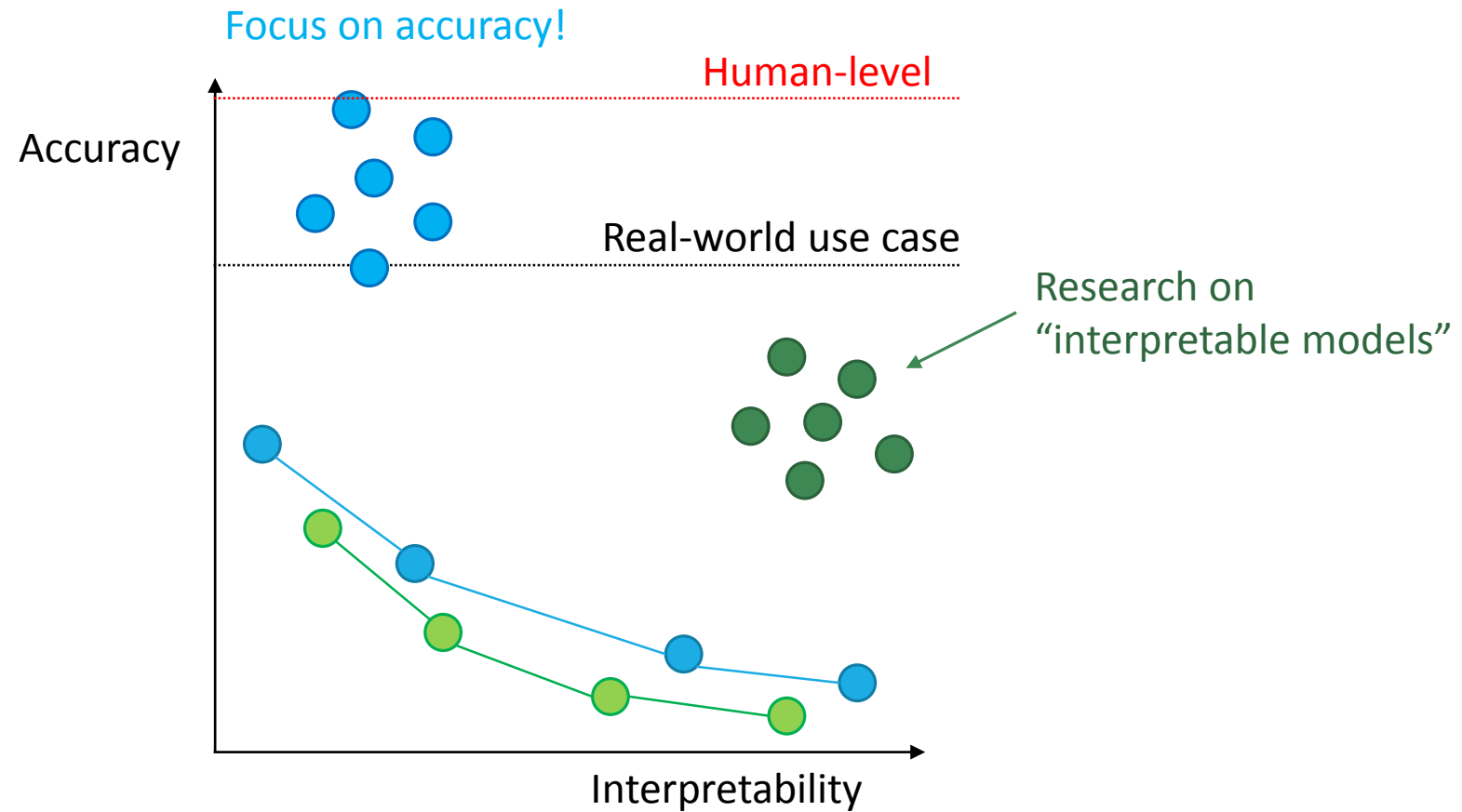
Lots of dimensions

**Black-boxes!**

# Accuracy vs Interpretability

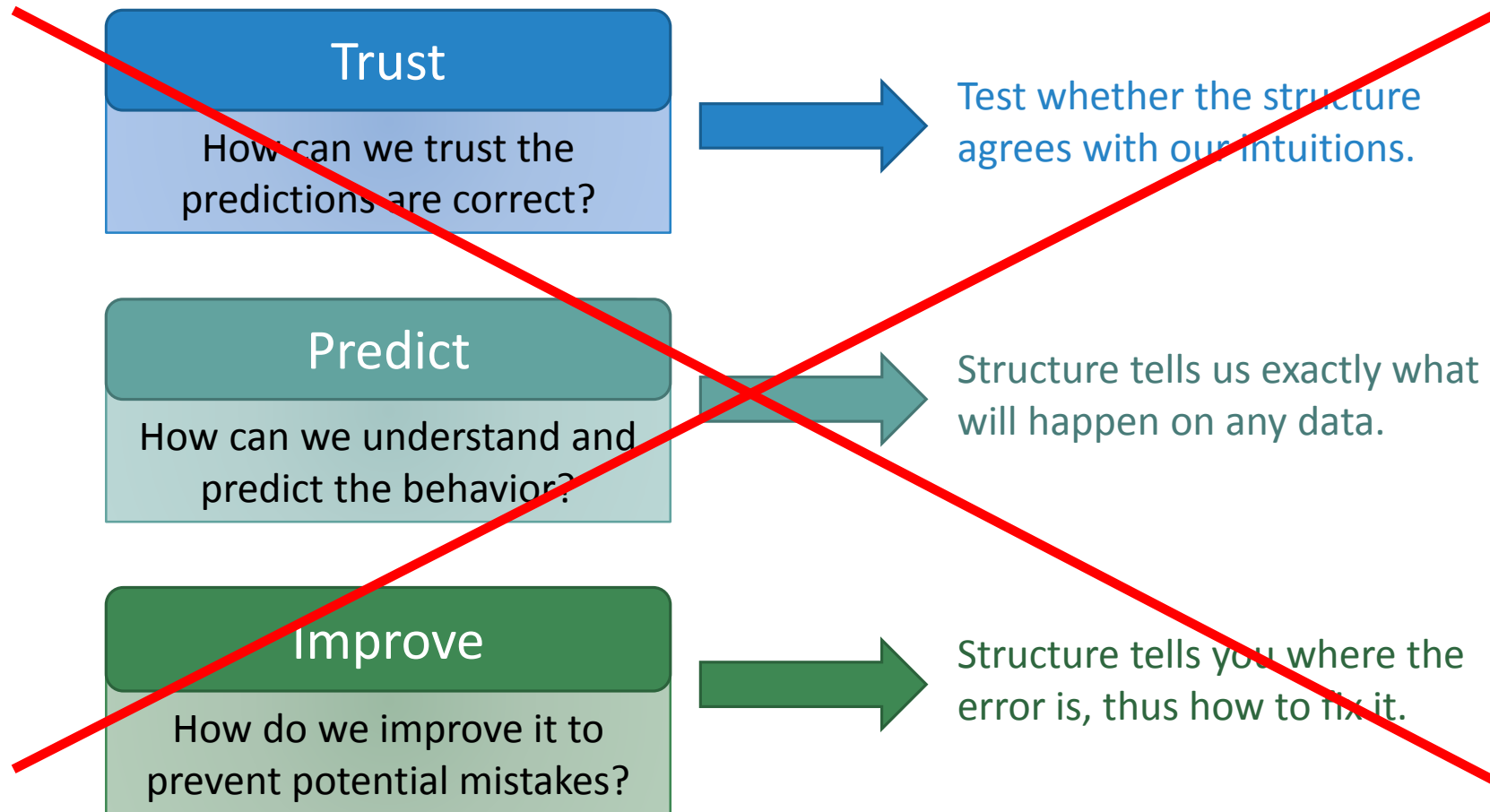


# Deep Learning



# Looking at the structure

---



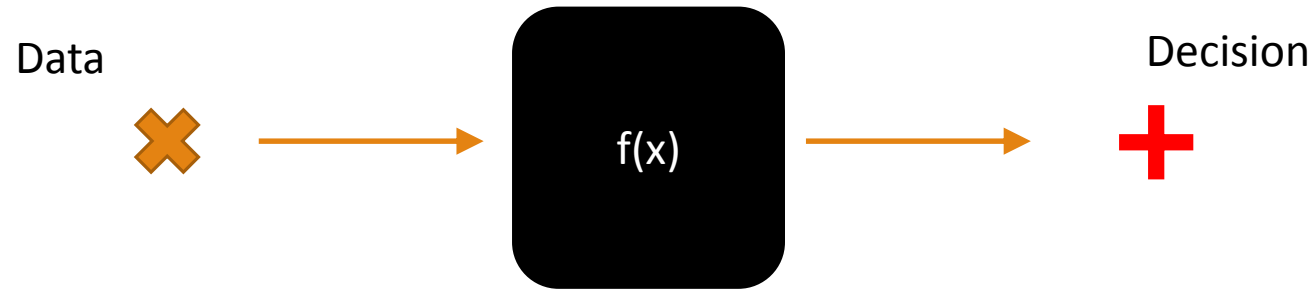
# Explaining Predictions

The LIME Algorithm

# Being Model-Agnostic...

---

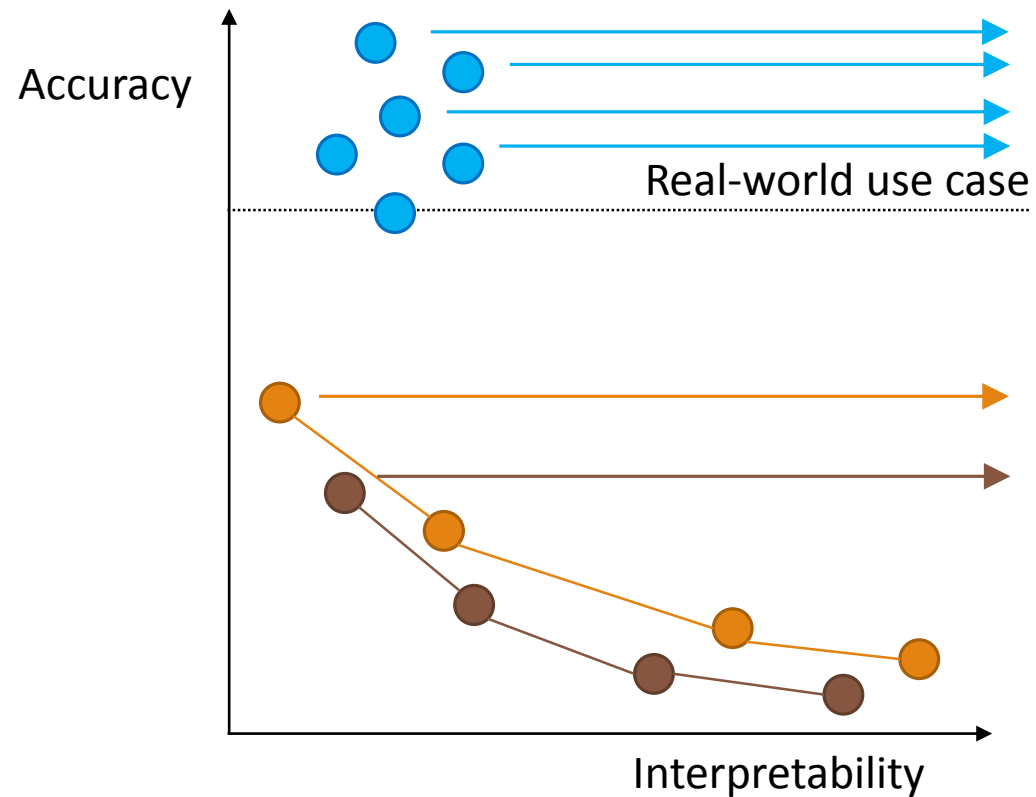
No assumptions about the internal structure...



Explain any existing, *or future*, model

# LIME: Explain Any Classifier!

---



Make  
everything  
interpretable!



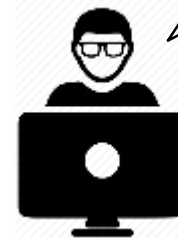
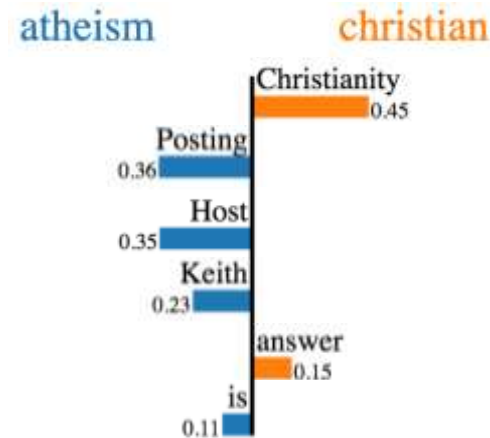
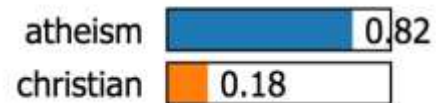
# What an explanation looks like

From: Keith Richards  
Subject: Christianity is the answer  
NTTP-Posting-Host: x.x.com

I think Christianity is the one true religion.  
If you'd like to know more, send me a note



Prediction probabilities

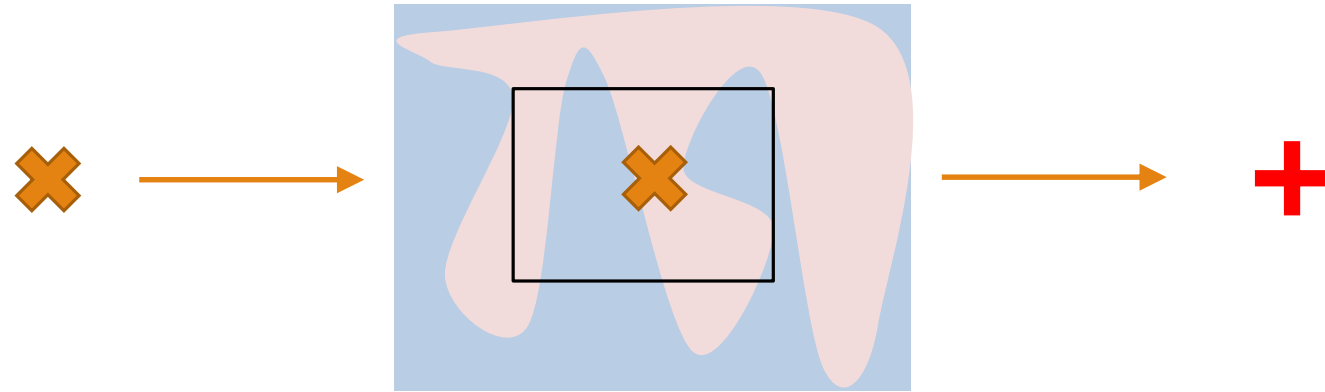


Why did this happen?

# Being Model-Agnostic...

---

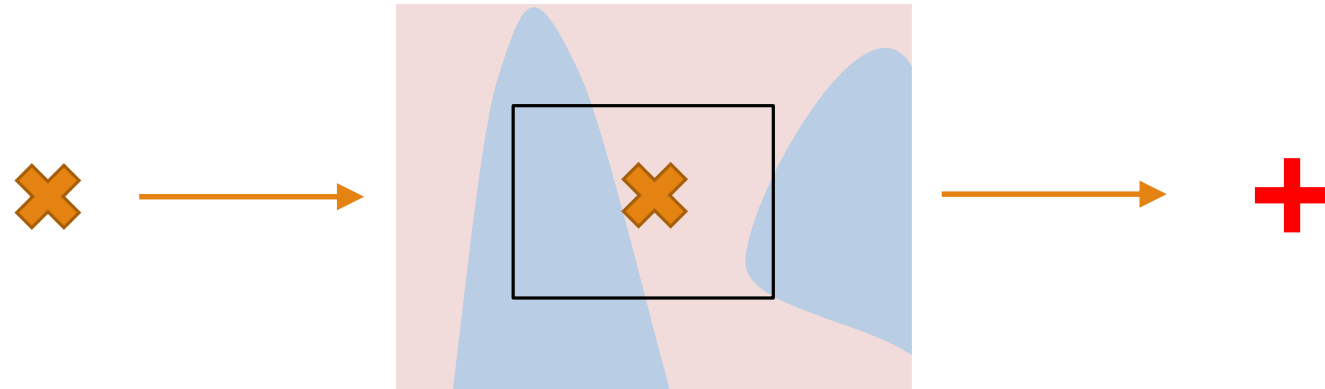
“Global” explanation is too complicated



# Being Model-Agnostic...

---

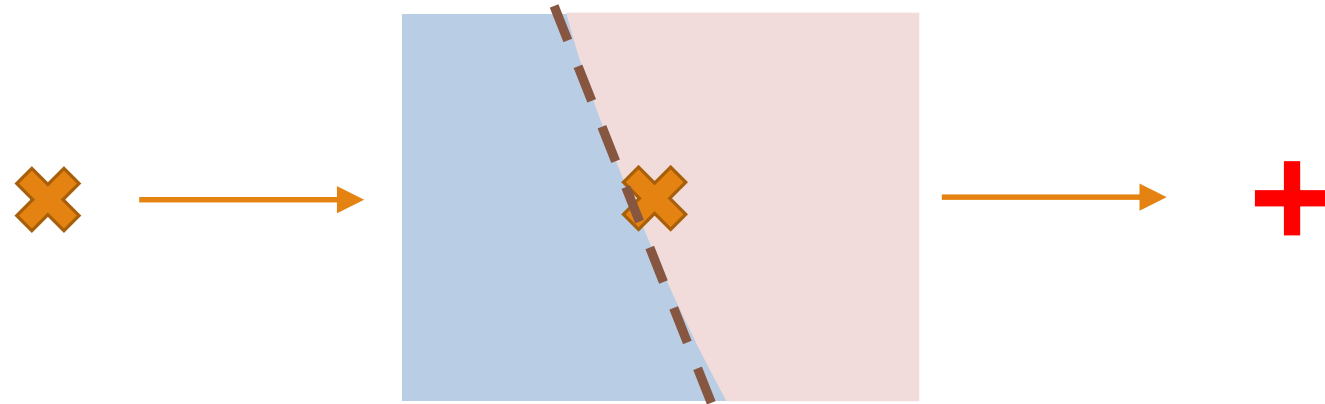
“Global” explanation is too complicated



# Being Model-Agnostic...

---

“Global” explanation is too complicated



Explanation is an interpretable model,  
that is locally accurate

# Google's Object Detector

---



$P(\text{🎸}) = 0.32$



$P(\text{🎸}) = 0.24$



$P(\text{🐶}) = 0.21$



# Classification: Wolf or a Husky?



Predicted: **wolf**  
True: **wolf**



Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**

Only 1 mistake!



Predicted: **wolf**  
True: **husky**



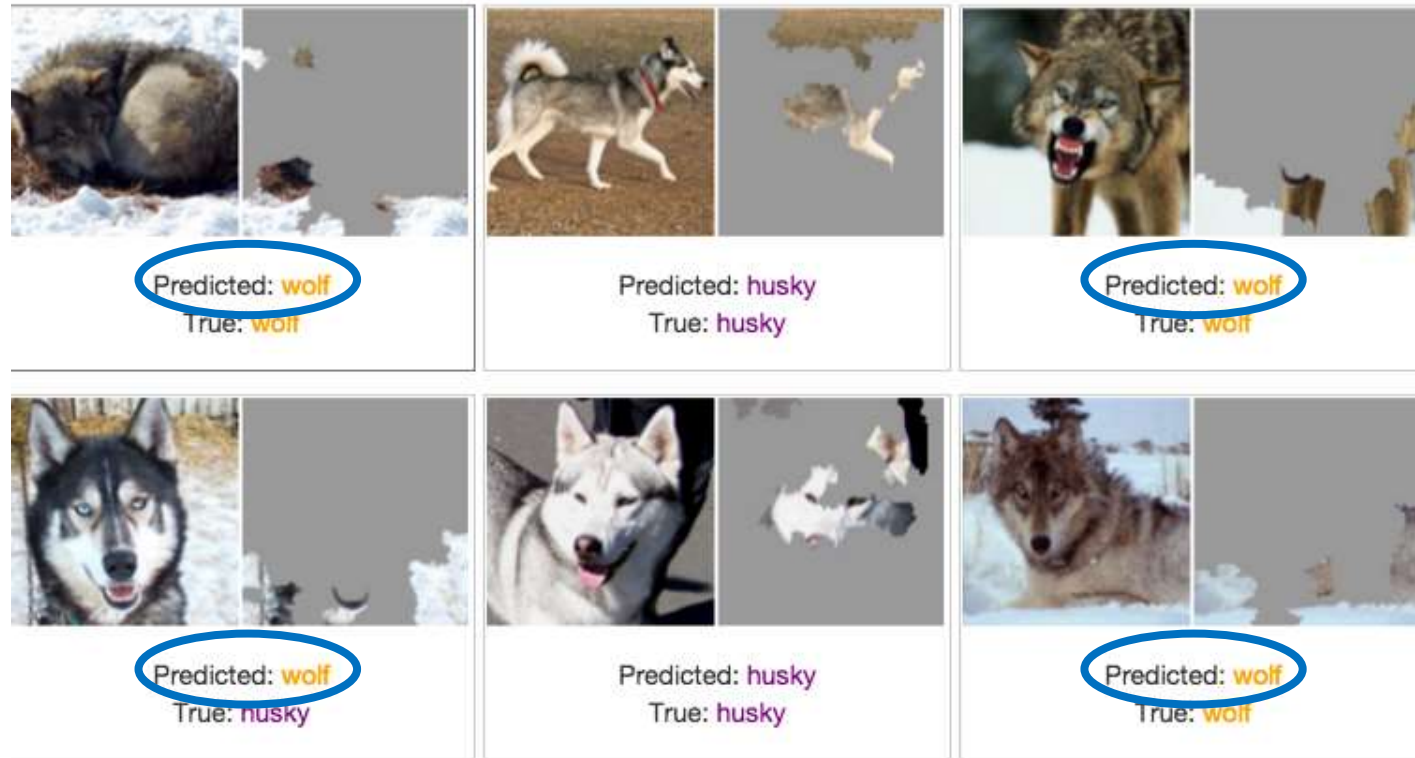
Predicted: **husky**  
True: **husky**



Predicted: **wolf**  
True: **wolf**



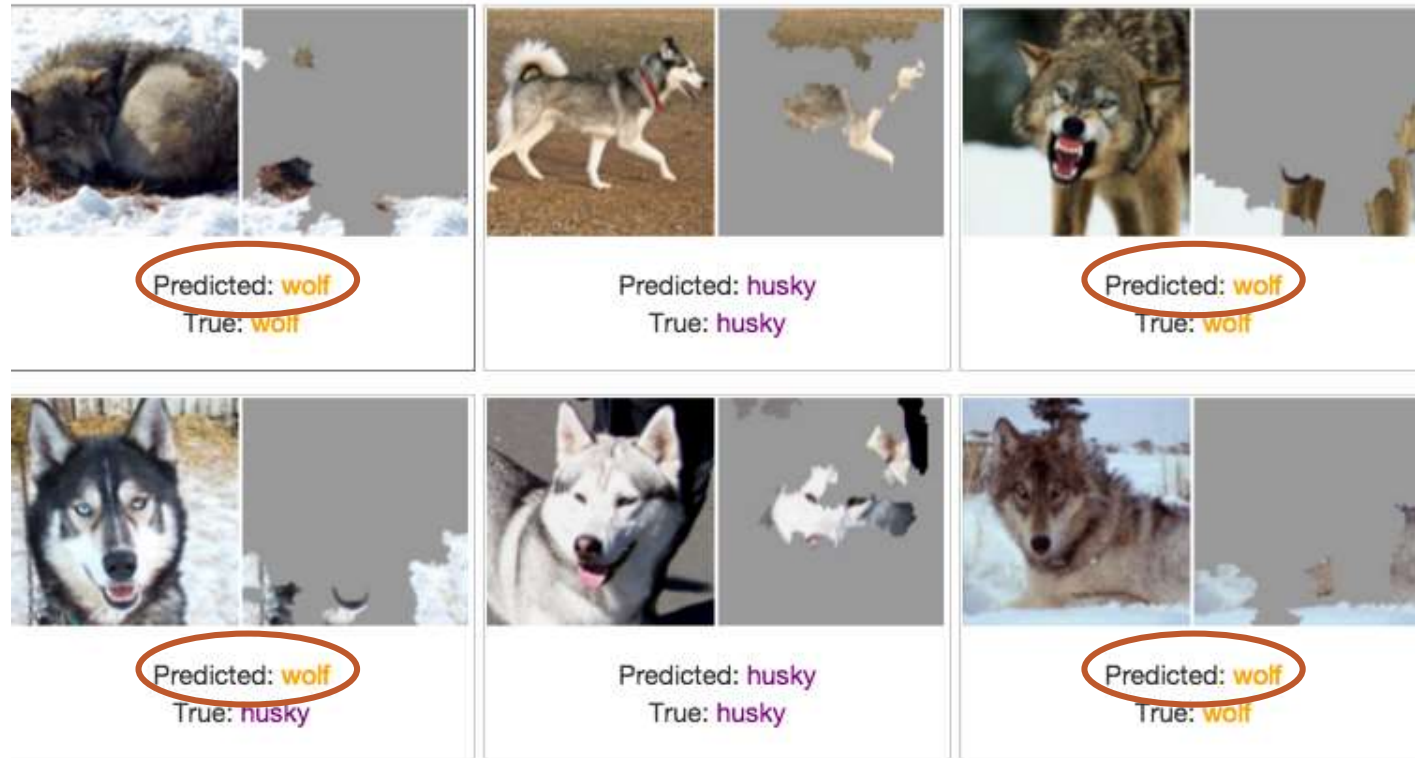
# Neural Network Explanations



We've built a great snow detector...



# Understanding Behavior



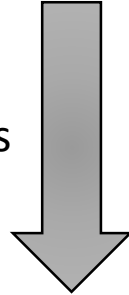
We've built a great snow detector...

# Comparing Classifiers

---

Classifier 1

Change the model  
Different data  
Different parameters  
Different “features”  
...



Classifier 2

Accuracy?

Look at Examples?

Deploy and Check?

“I have a gut feeling..”

Explanations?

# Comparing Classifiers

---



Original Image



"Bad" Classifier



"Good" Classifier

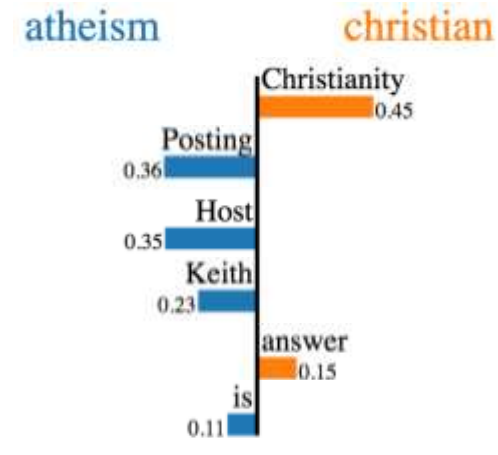
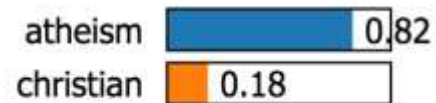
# Explanation for a bad classifier

From: Keith Richards  
Subject: Christianity is the answer  
NTTP-Posting-Host: x.x.com

I think Christianity is the one true religion.  
If you'd like to know more, send me a note

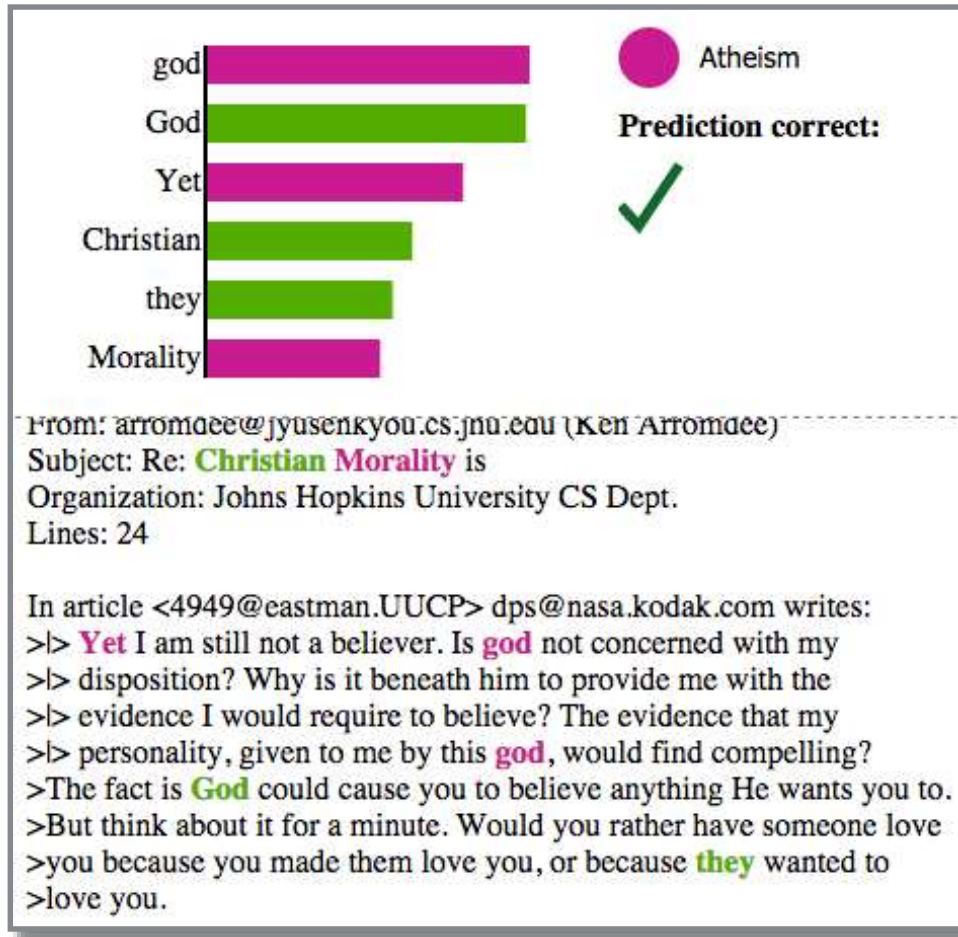


Prediction probabilities



After looking at the explanation,  
we shouldn't trust the model!

# “Good” Explanation



It seems to be picking up on more reasonable things.. good!

# Recent Work

Counter-examples and Counter-factuals

# Understanding via Predicting

---

Users “understand” a model if they can predict its behavior on unseen instances

Precision

How accurate are the users guesses?  
If the users guess wrong, they don't understand

Coverage

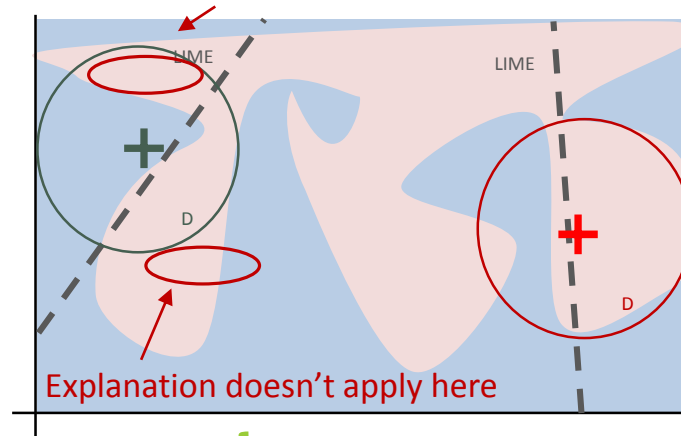
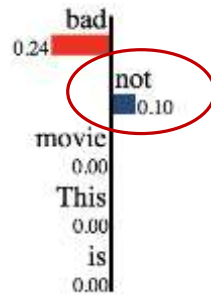
How often do the users make confident guesses?  
It's okay not to be able to guess!

Precision is much more  
important than Coverage!

It's much better not to guess than to guess  
confidently, but be completely wrong!

# Linear Explanations

+ This movie is not bad. Explanation is wrong in this region

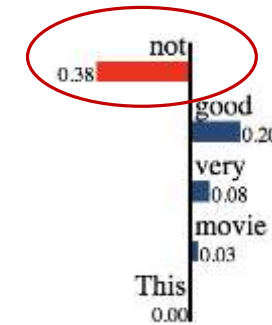


Explanation doesn't apply here

D

This director is always bad.  
This movie is not nice.  
This stuff is rather honest.  
This star is not bad.  
...

+ This movie is not very good.



This explanation is a better approximation than the other one.

**Problem 1:** Where is the explanation good?

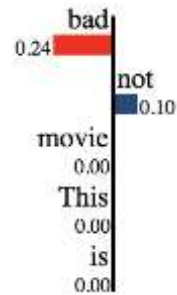
**Problem 2:** What is the coverage?

→ Users will make mistakes!



# Anchors: Precise Counter-factuals

+ This movie is not bad.

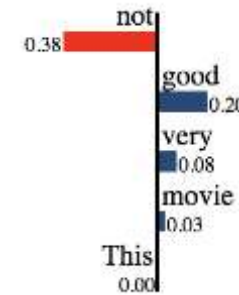


Anchor: "not bad" → Positive

$D(\cdot|A)$  { This audio is not bad.  
This novel is not bad.  
This footage is not bad.

An anchor is a sufficient condition

+ This movie is not very good.



Anchor: "not good" → Negative

$D(\cdot|A)$  { This poster is not ever good.  
This picture is not rarely good.  
This actor is not incredibly good.

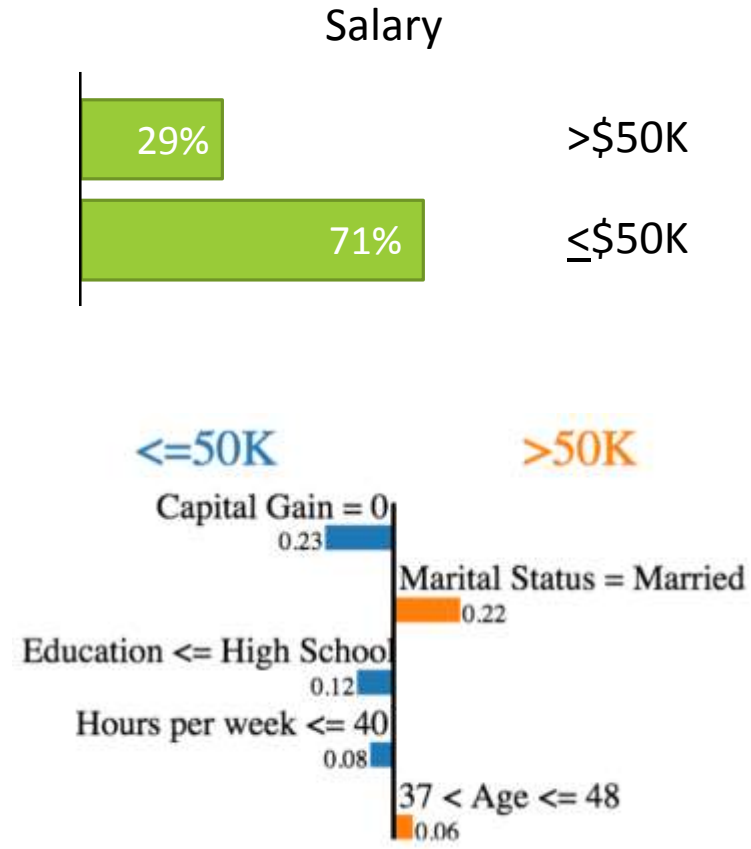
Clear (and adaptive) coverage

Probabilistic guarantee avoids human mistakes

# Salary Prediction

Feature	Value
Age	$37 < \text{Age} \leq 48$
Workclass	Private
Education	$\leq$ High School
Marital Status	Married
Occupation	Craft-repair
Relationship	Husband
Race	Black
Sex	Male
Capital Gain	0
Capital Loss	0
Hours per week	$\leq 40$
Country	United States

IF Education  $\leq$  High School  
Then Predict Salary  $\leq$  50K



# Visual QA

---



---

**What** is the mustache made of?    banana

---

---

How **many** bananas are in the picture?    2

---

# Encoder/Decoder LSTMs

---

English	Portuguese
<b>This</b> is the <b>question</b> we must address	<b>Esta</b> é a questão que temos que enfrentar.

# Encoder/Decoder LSTMs

---

English	Portuguese
<b>This is the question</b> we must address	<b>Esta</b> é a questão que temos que enfrentar.
<b>This is the problem</b> we must address	<b>Este</b> é o problema que temos que enfrentar.

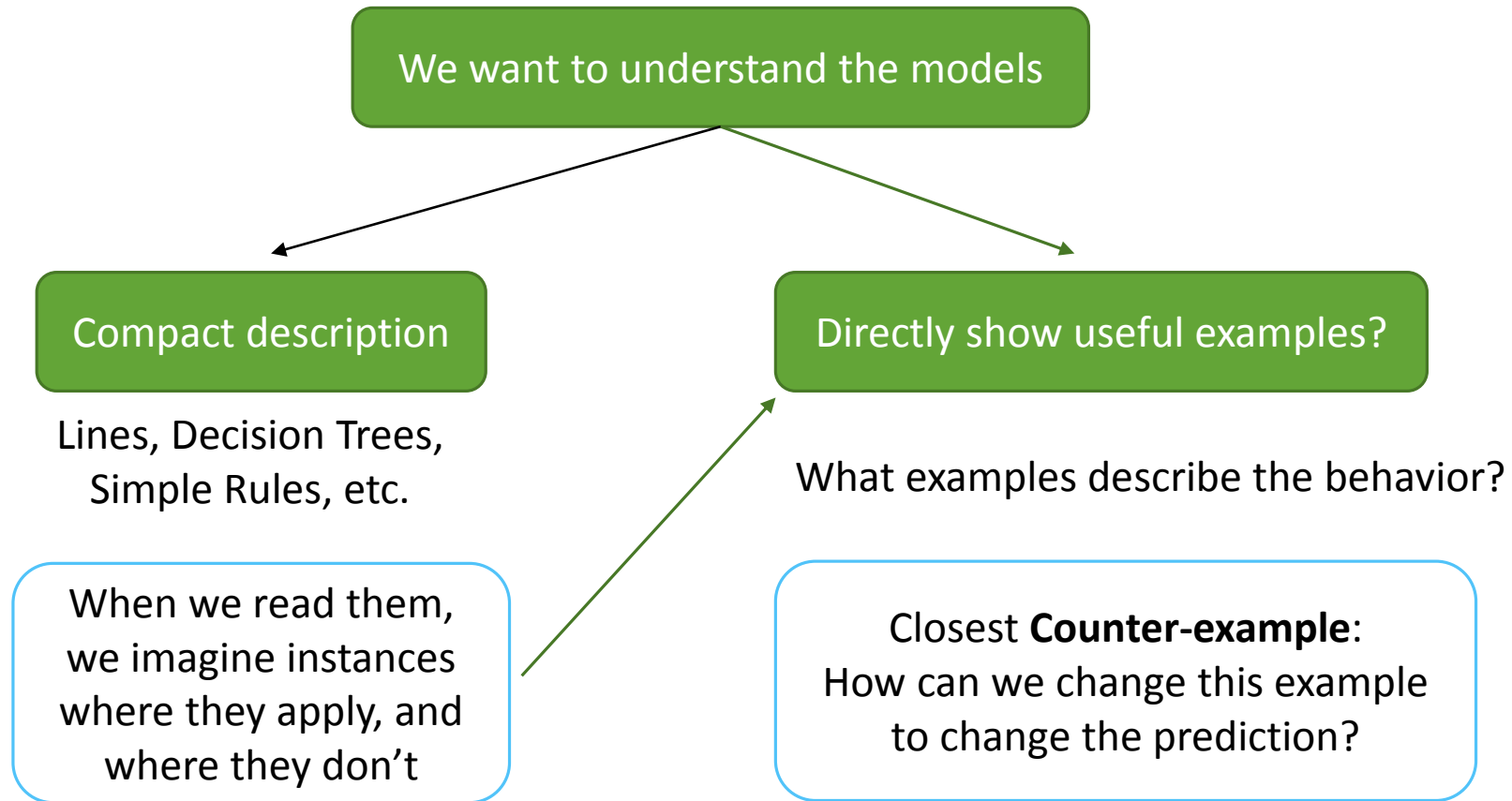
# Encoder/Decoder LSTMs

---

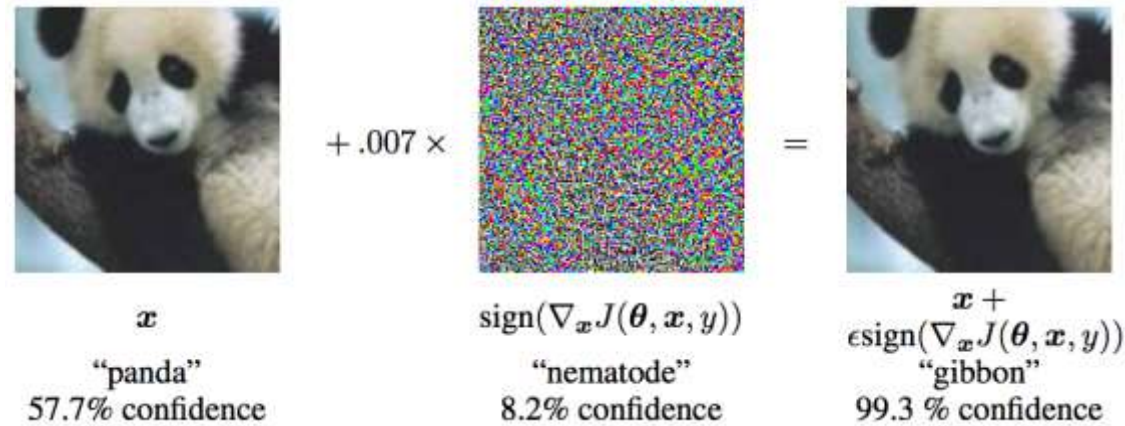
English	Portuguese
<b>This is the question</b> we must address	<b>Esta</b> é a questão que temos que enfrentar.
<b>This is the problem</b> we must address	<b>Este</b> é o problema que temos que enfrentar.
<b>This is what</b> we must address	É <b>isso</b> que temos de enfrentar.

# What's a Good Explanation?

---

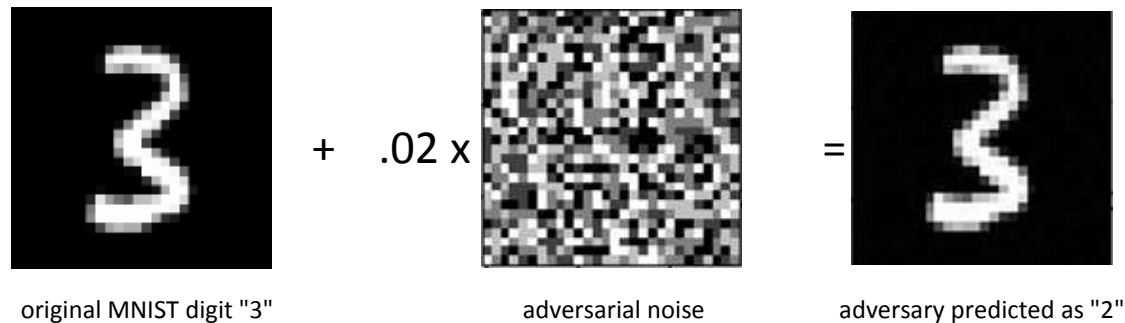


# Adversarial Examples


$$\begin{array}{ccc} \text{panda} & + .007 \times & \text{nematode} & = & \text{gibbon} \\ \text{57.7\% confidence} & & \text{8.2\% confidence} & & \text{99.3\% confidence} \\ x & & \text{sign}(\nabla_x J(\theta, x, y)) & & x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)) \end{array}$$

"inputs formed by applying small but intentionally worst-case perturbations to examples from the dataset, such that the perturbed input results in the model outputting an incorrect answer with high confidence"

Goodfellow et al, "Explaining and Harnessing Adversarial Examples", ICLR 2015.


$$\begin{array}{ccc} \text{original MNIST digit "3"} & + .02 \times & \text{adversarial noise} & = & \text{adversary predicted as "2"} \end{array}$$



# Adversarial Examples: Pros

---

$$x^* = \operatorname{argmin}_{\tilde{x}} \|x - \tilde{x}\|_2 \text{ s.t. } f(x) \neq f(\tilde{x})$$

## Advantages:

- Applicable to any gradient -based classifier
- Useful to evaluate the *robustness* of the model against adversaries
- Small perturbations often lead to imperceivable adversarial examples

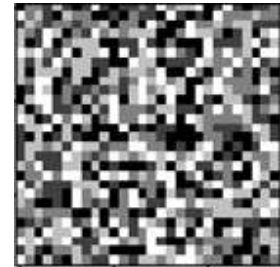
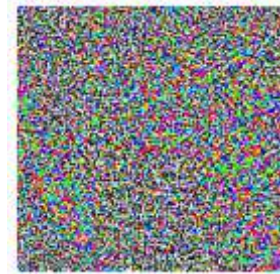
# Adversarial Examples: Cons

---

$$x^* = \operatorname{argmin}_{\tilde{x}} \|x - \tilde{x}\|_2 \text{ s.t. } f(x) \neq f(\tilde{x})$$

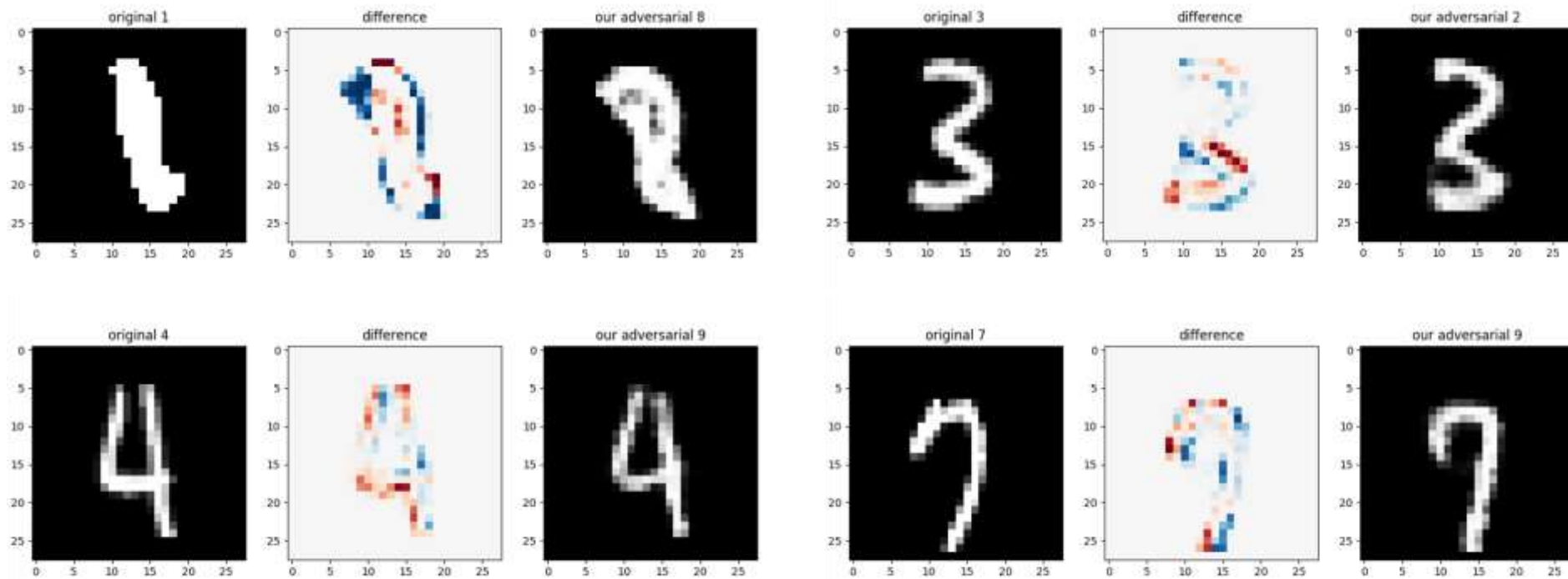
## Disadvantages:

- Examples are unnatural
  - may not look anything you would naturally see in the "wild"
- Distance is not always meaningful
  - E.g. color change or translation/rotation of an image
- Cannot be used for structured domains like text, code, etc.:
  - E.g. replacing/removing words results in sentences that are not grammatical
- Do not provide insights into why the sample is an adversary
  - How is the model working?
  - How to fix the model?



# Example: MNIST Digits

---



# Example: Church vs Tower

church→tower



tower→church



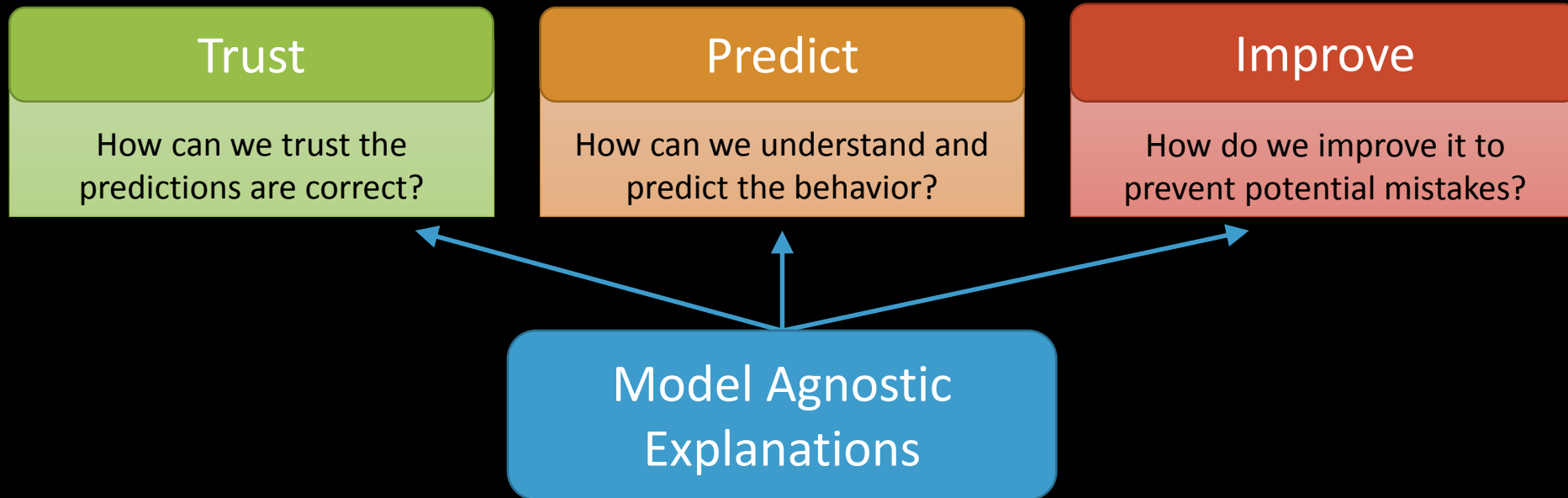
# Machine Translation

---

Debug Google Translate, remotely!

Source Sentence (English)	Generated Translation (German)
s : People sitting in a dim restaurant <b>eating</b>	Leute, die in einem dim Restaurant <b>essen</b> sitzen.
s' : People sitting in a living room <b>eating</b> .	Leute, die in einem Wohnzimmeressen sitzen. <i>(People sitting in a living room)</i>
s : Elderly people <b>walking</b> down a city street .	Ältere Menschen, die eine Stadtstraße <b>hinuntergehen</b> .
s' : A man <b>walking</b> down a street playing	Ein Mann, der eine Straße entlang spielt. <i>(A man playing along a street.)</i>

# Explanations are important!



## Model Agnostic Explanations

Work with Marco T. Ribeiro, Carlos Guestrin, Dheeru Dua, and Zhengli Zhao

# Thanks!

sameer@uci.edu  
sameersingh.org