

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

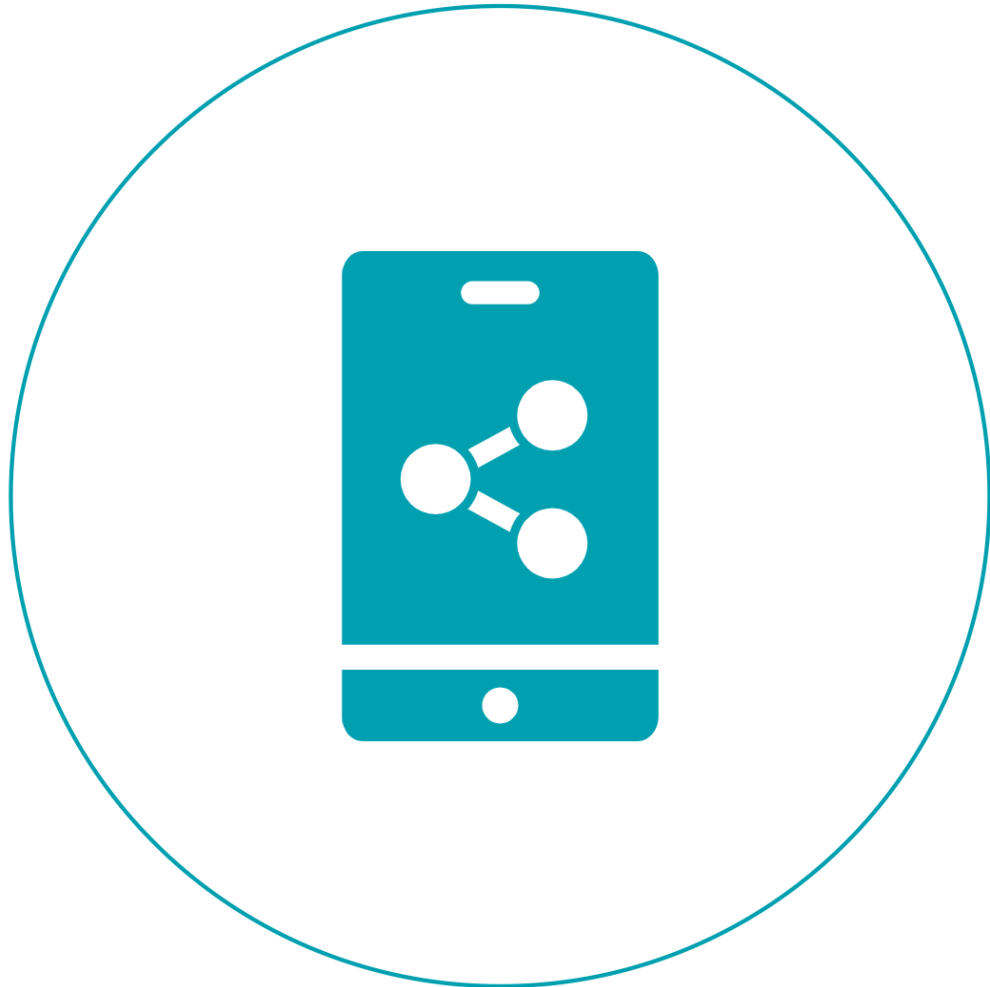
Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

Fengmin Wu

2019.10.16

Why is interpretability discussed today?

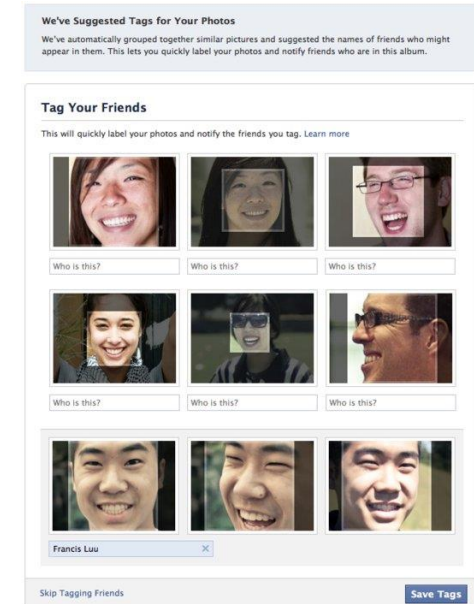
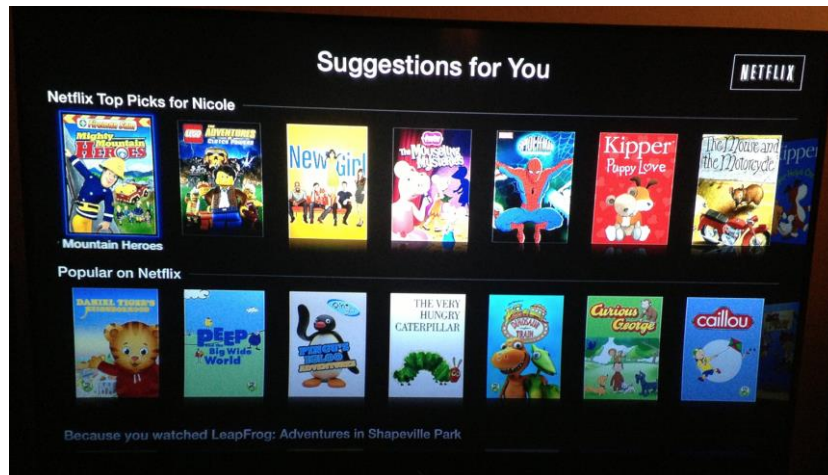
... since AI and machine learning is not really new ...



AI and machine learning algorithms become omnipresent.

- In (almost) every area of our lives
- Because algorithms tend to be more efficient and faster than humans

Machine Learning applications becoming pervasive...



Why is interpretability discussed today?

... since AI and machine learning is not really new ...

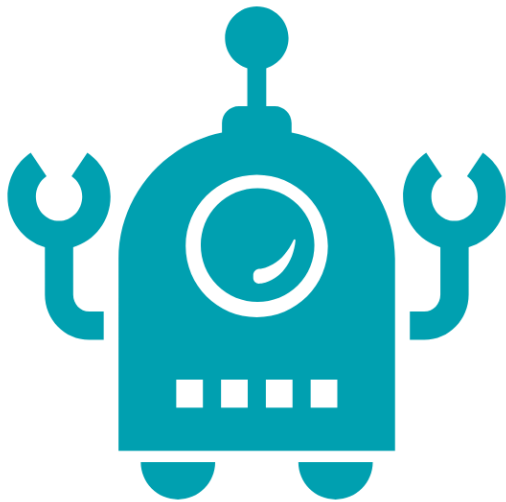


Data privacy and security become more important.

- Majority of our society is touched by it
- Talked about in the media
- Thus became an important political topic

Why is interpretability discussed today?

... since AI and machine learning is not really new ...

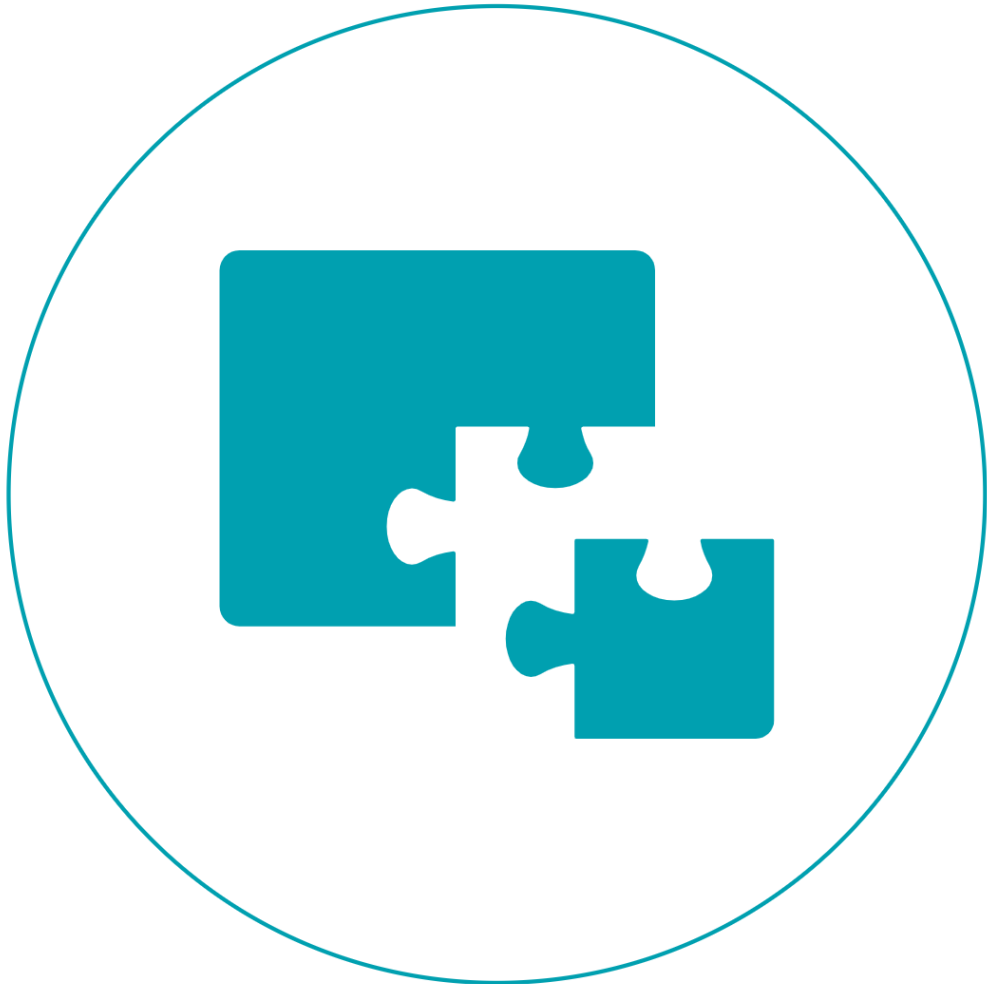


AI is hyped.

- The media reports about it a lot
- And tends to be critical about it

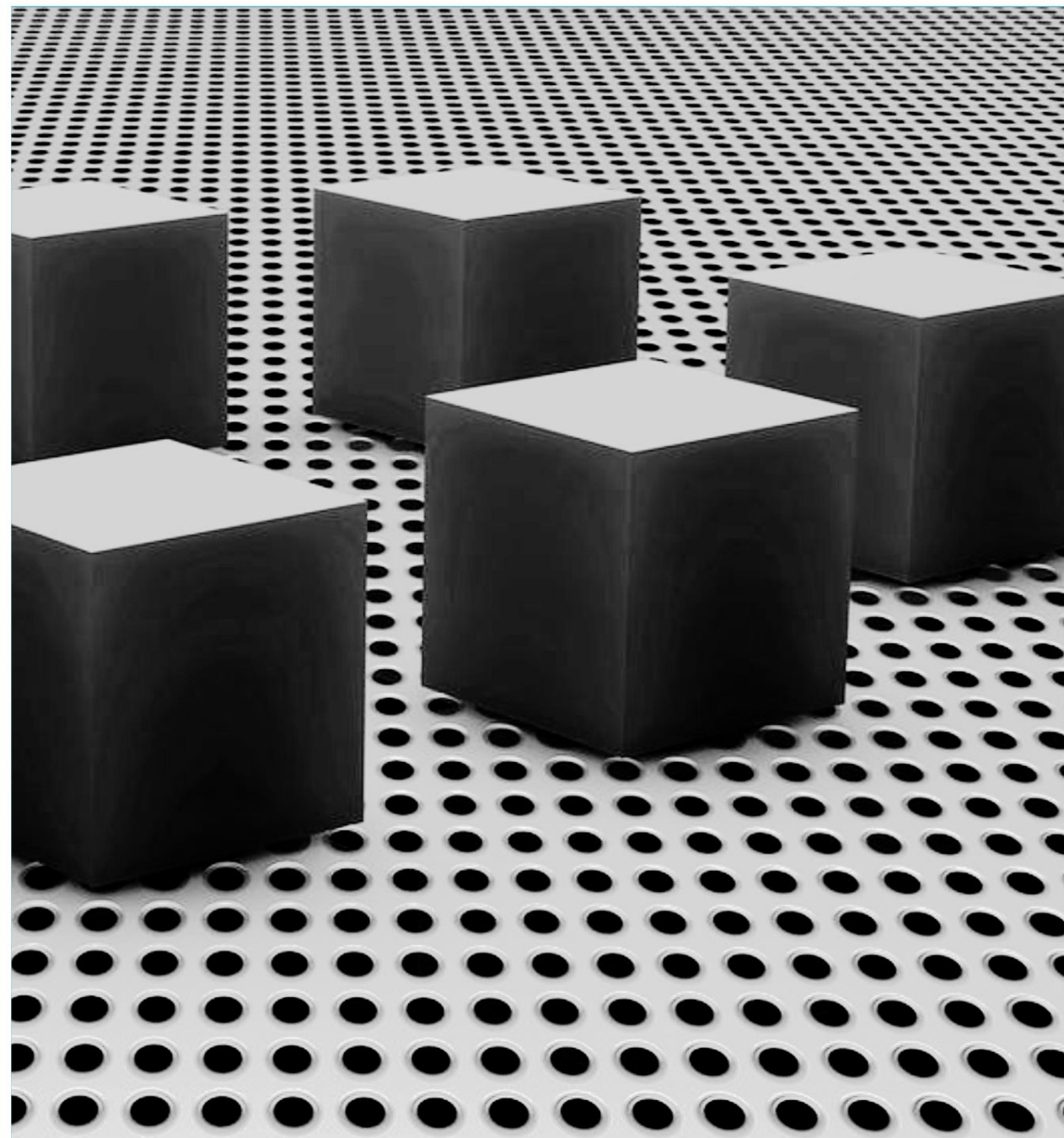
Why is interpretability discussed today?

... since AI and machine learning is not really new ...



AI is complex.

- Possibilities and dangers are hard to differentiate by most people
- “Brave New World”?

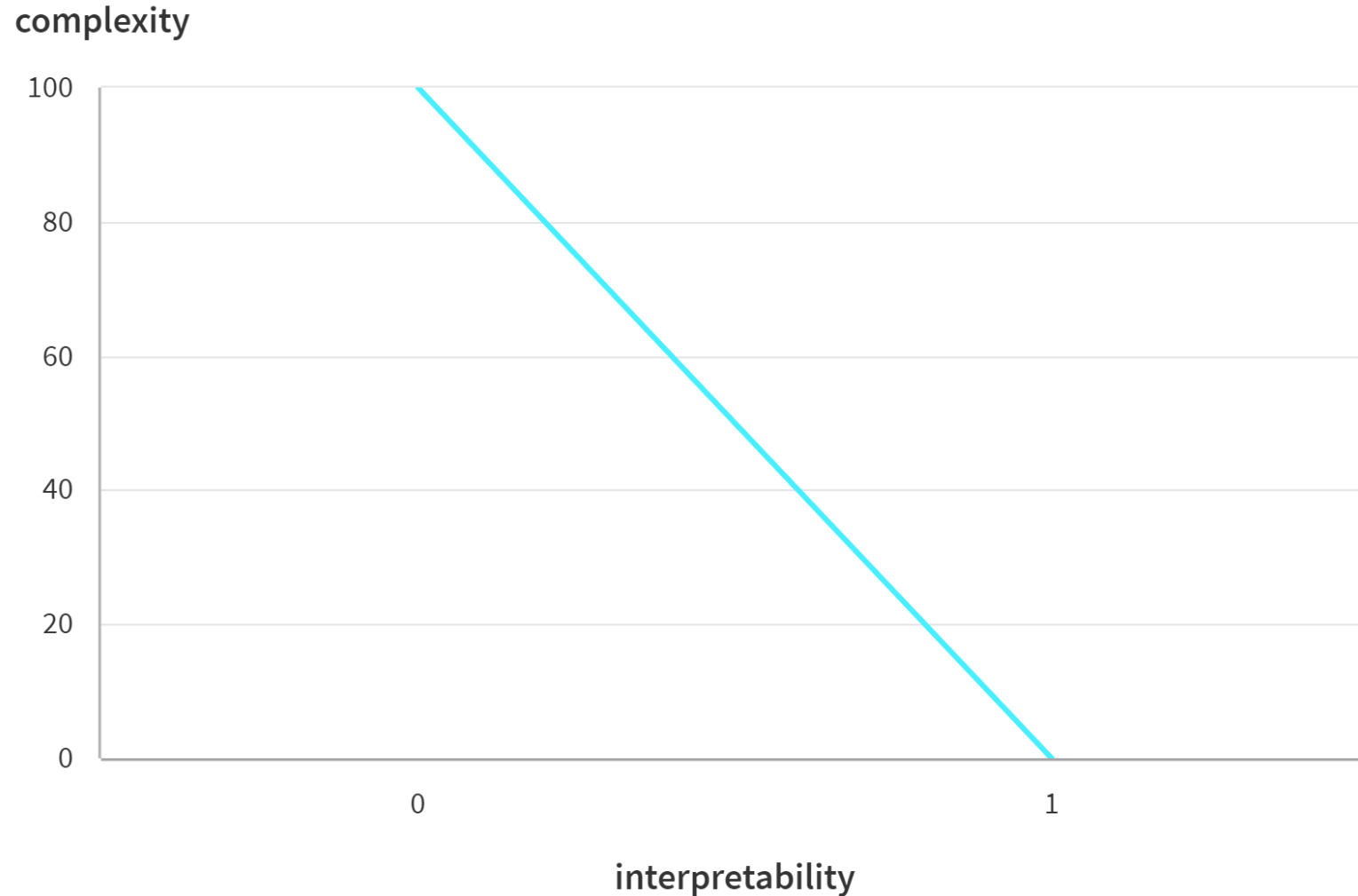


Why are ML models

“Black-Boxes”
?

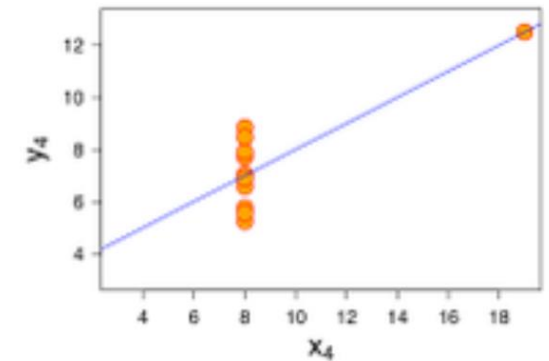
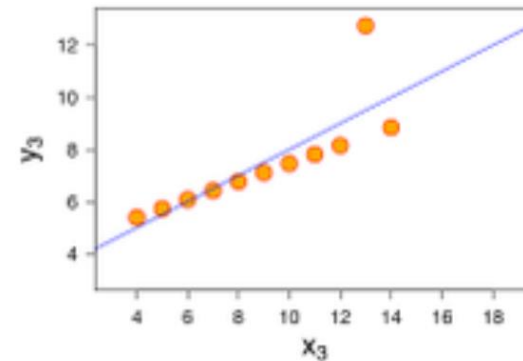
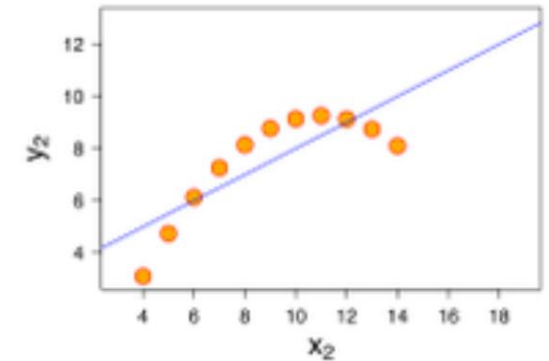
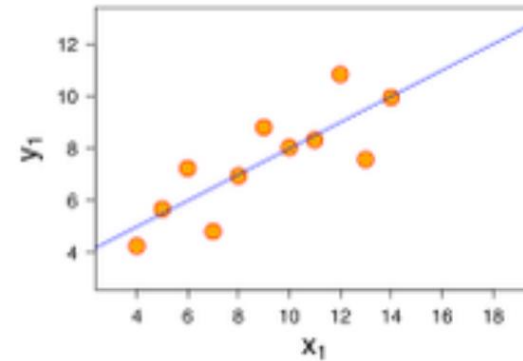
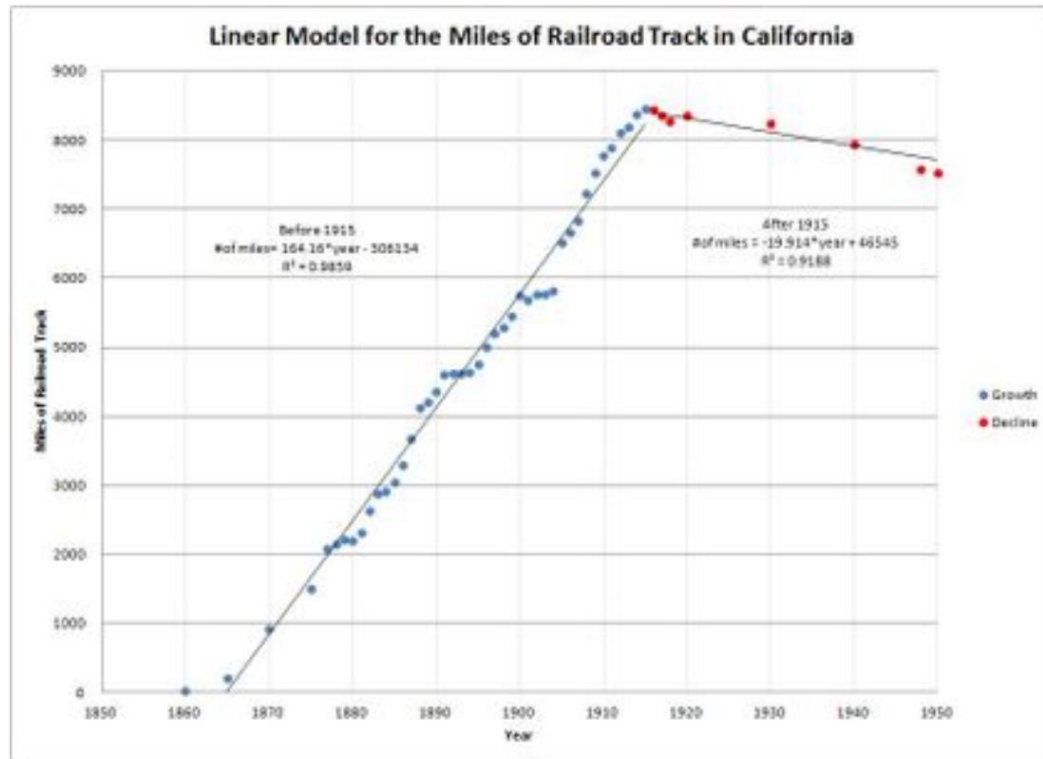
Trade-off between interpretability & complexity

The more complex a model, the harder it tends to be to understand.



Trade-off between interpretability & complexity

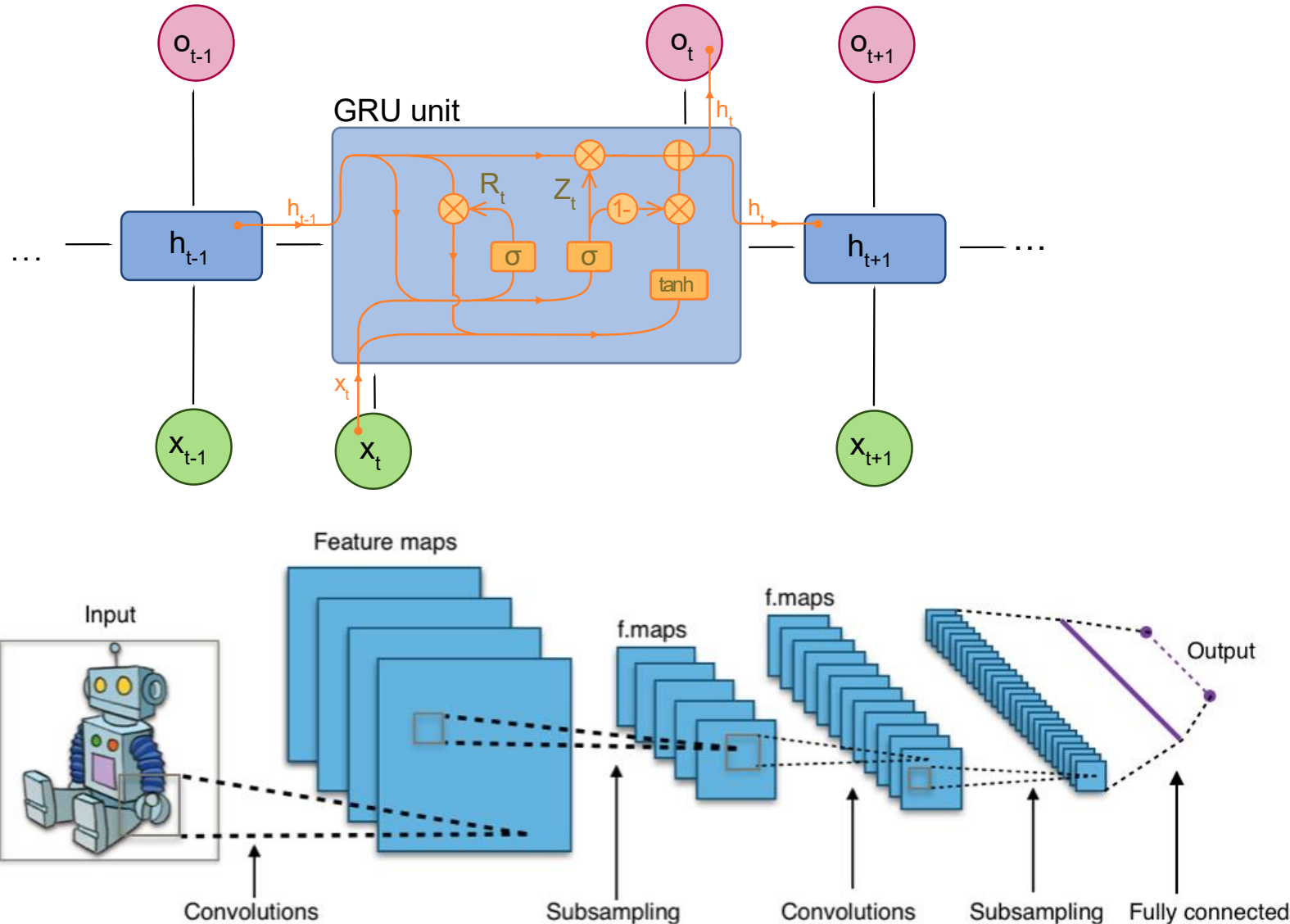
Linear regression models are easy to understand: If ..., then ...



Source: Wikipedia

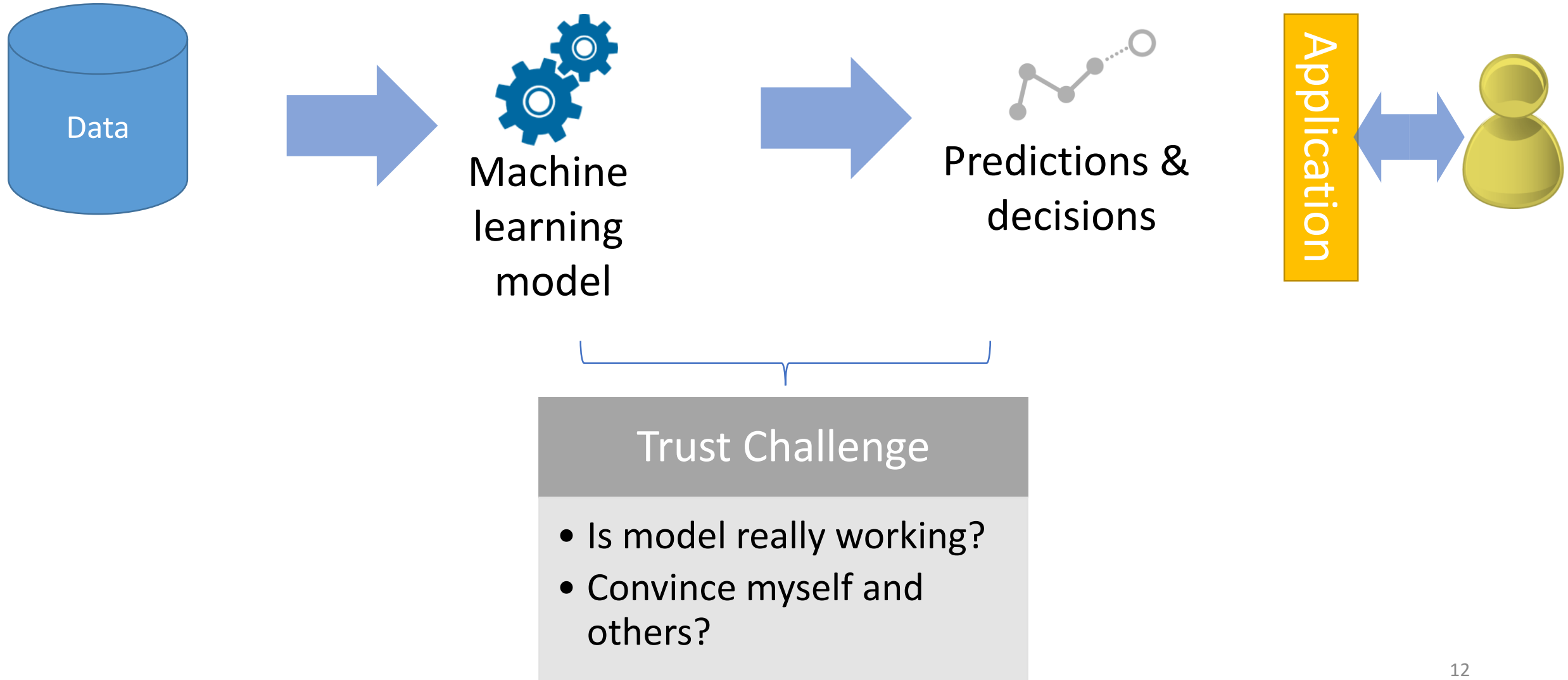
Trade-off between interpretability & complexity

Complexity tends to make models perform better but harder to understand.



Can we **trust** our models?
Are they fit for the wild?

How to build an application with ML



Why should we improve our understanding of ML models?

... if technically it isn't necessary ...



Improving our models

Generalisability

“Sanity Check”

Prevent wrong conclusions &
potentially adversarial attacks

Improving our models

When we understand our models better, we are better at detecting wrong conclusions.



Example 1: Image classification of wolves & Huskies

The model based its predictions on the snow in the background.



Example 2: Text classification of Christian & Atheist posts

Not all words that predictions were based on made sense.



Example 3: Image classification with Google's Inception Net

Knowing which areas of an image contributed to a decision helps us trust it.

Why should we improve our understanding of ML models?

... if technically it isn't necessary ...



Improving our models

Generalisability

“Sanity Check”

Prevent wrong conclusions &
potentially adversarial attacks



Trust and transparency

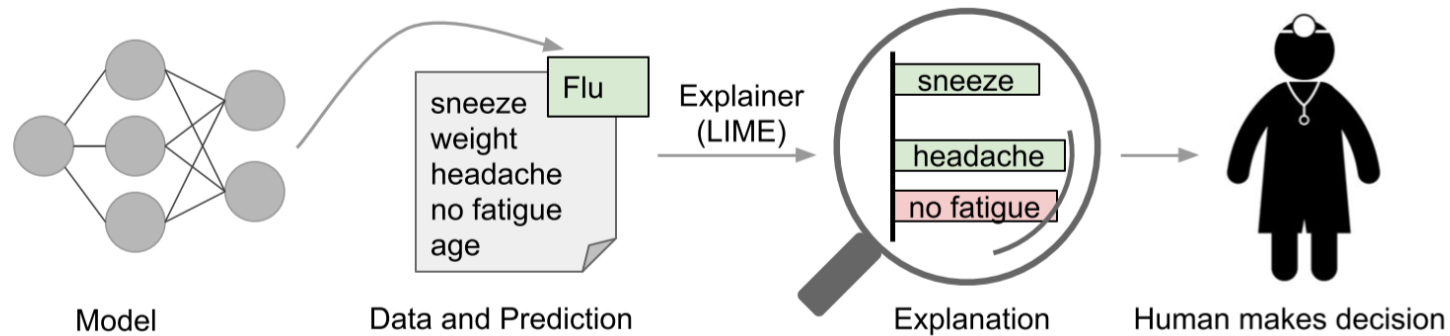
Can I trust my model's decisions?

Why does my model make the predictions it
makes?

Trust and transparency

But trust is also important in others areas!

- Decision can be a question of life and death.
- Medical intervention needs to be based on a diagnosis.



Trust and transparency

But trust is also important in others areas!

- Decision can be a question of life and death.
- Medical intervention needs to be based on a diagnosis.
- In business, we can save time and money by improving our understanding of machine learning.



Why should we improve our understanding of ML models?

... if technically it isn't necessary ...



Improving our models

Generalisability

“Sanity Check”

Prevent wrong conclusions &
potentially adversarial attacks



Trust and transparency

Can I trust my model's decisions?

Why does my model make the predictions
it makes?



Prevent Bias

Fairness

Identify and prevent bias

Explaining individual predictions:

Making any model interpretable

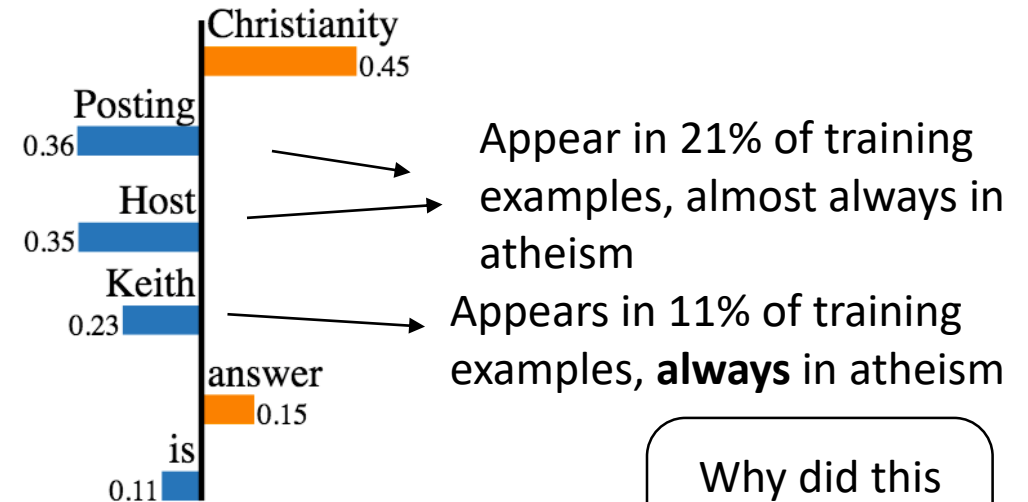
What an explanation looks like

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

atheism

christian



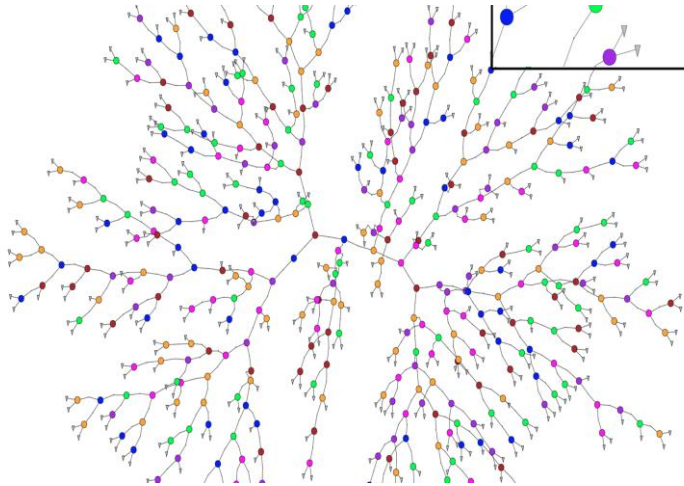
Why did this happen? How do I fix it?

→ Will not generalize
→ Don't trust this model!

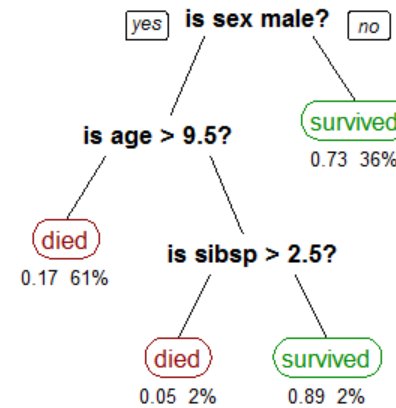
Three must-haves for a good explanation

Interpretable

- Humans can easily interpret reasoning



Definitely
not interpretable



Potentially
interpretable

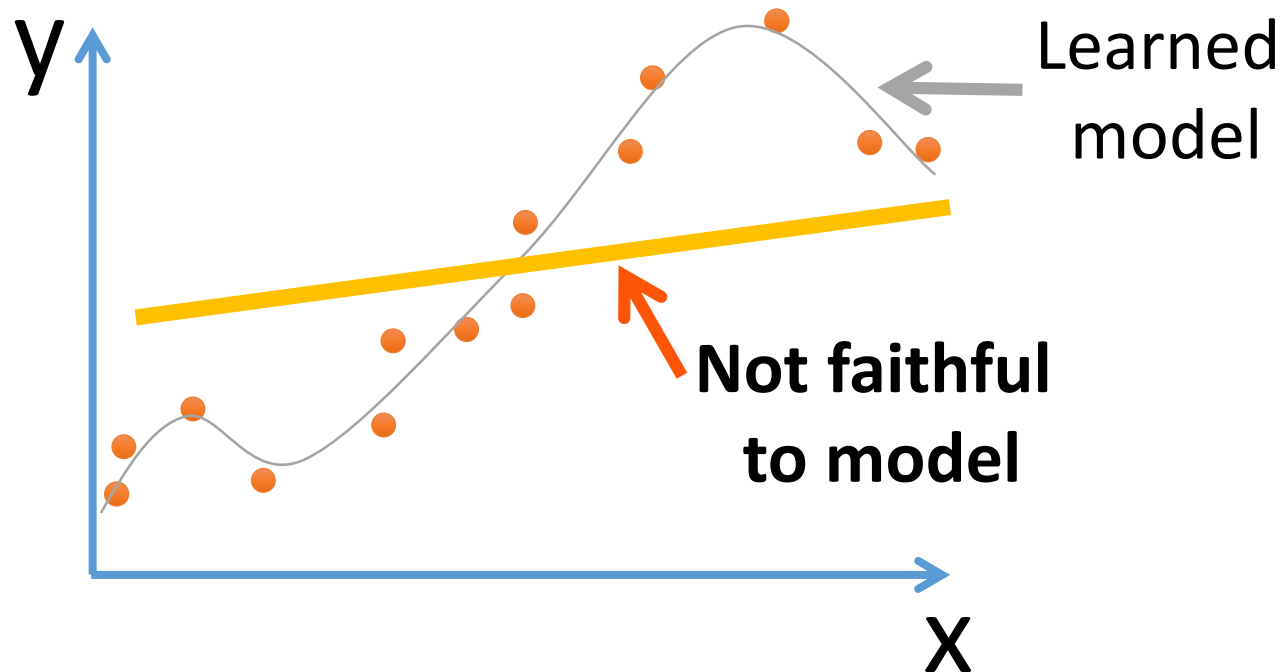
Three must-haves for a good explanation

Interpretable

- Humans can easily interpret reasoning

Faithful

- Describes how this model actually behaves



Three must-haves for a good explanation

Interpretable

- Humans can easily interpret reasoning

Faithful

- Describes how this model actually behaves

Model agnostic

- Can be used for *any* ML model

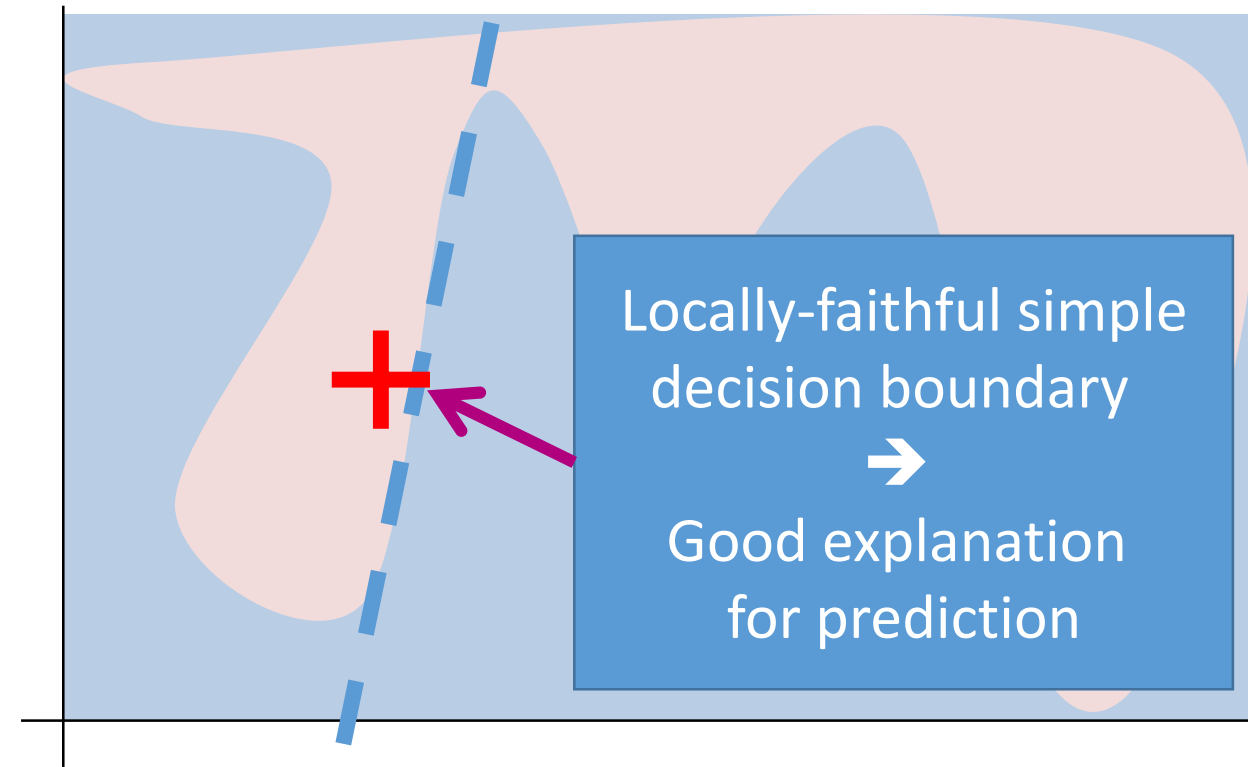


LIME: Local Interpretable Model-Agnostic Explanations

LIME – Key Ideas

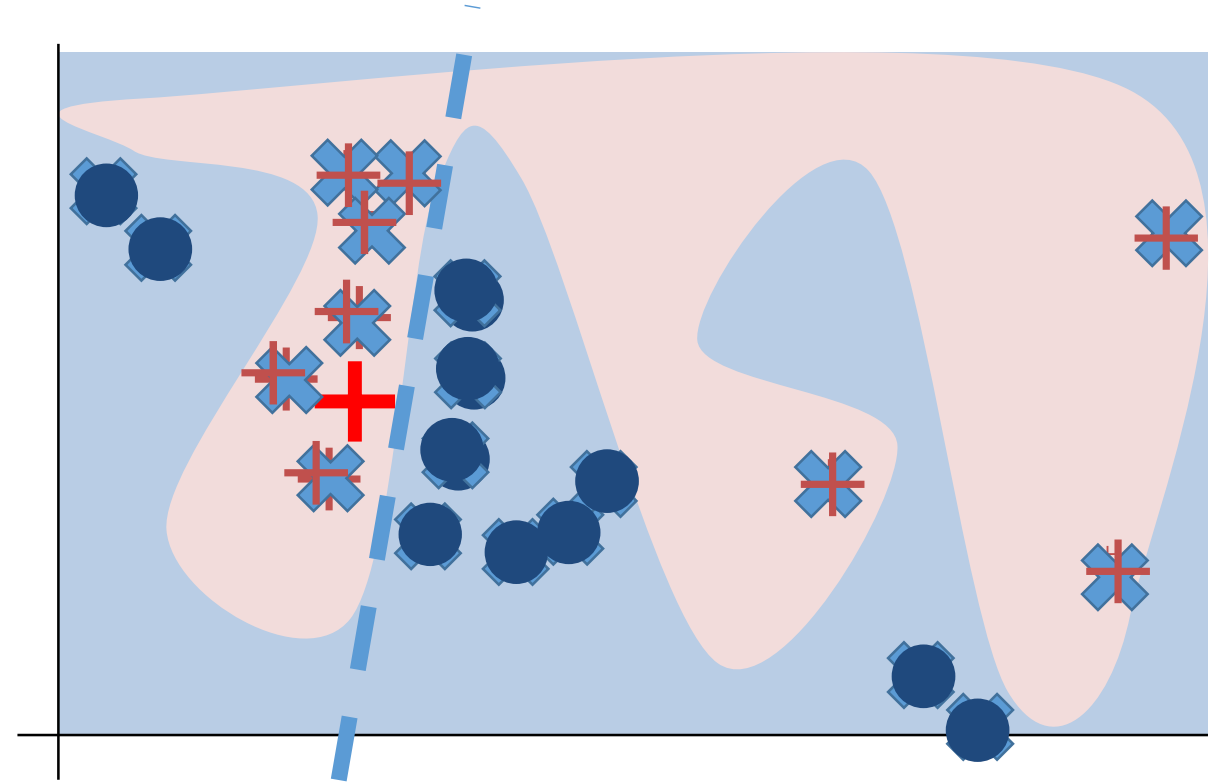
1. Pick a model class interpretable by humans
 - Not globally faithful... ☹️
2. Locally approximate global (blackbox) model
 - Simple model globally bad, but locally good

Line,
shallow decision tree,
sparse features, ...



Using LIME to explain a complex model's prediction for input x_i

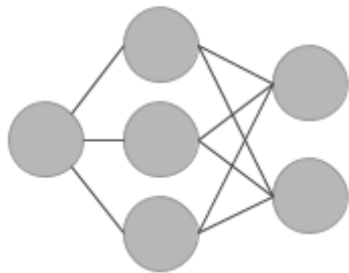
1. Sample points around x_i
2. Use complex model to predict labels for each sample
3. Weigh samples according to distance to x_i
4. Learn new simple model on weighted samples
5. Use simple model to explain



Interpretable representations

x (embeddings)

0.5	0.3	1.3	4.4	1.1	...
-----	-----	-----	-----	-----	-----



Model

This is what we perturb, and this is what we use in the interpretable approximation

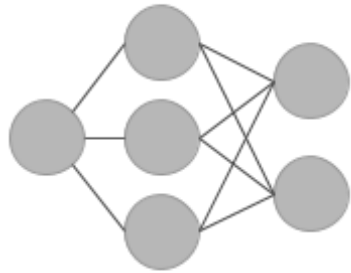
x' (words)

This is a horrible movie.



Interpretable representation: images

x (3 color channels / pixel)



Model

x' (contiguous superpixels)



Human

LIME



Explain model : $g \in G$

Model being explained : $f: \mathbb{R}^d \rightarrow \mathbb{R}$

A measure of complexity : $\Omega(g)$ - depth of the tree, number of non-zero weights




Capturing the locality : π_x - far away from x , low weight, or vice versa

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

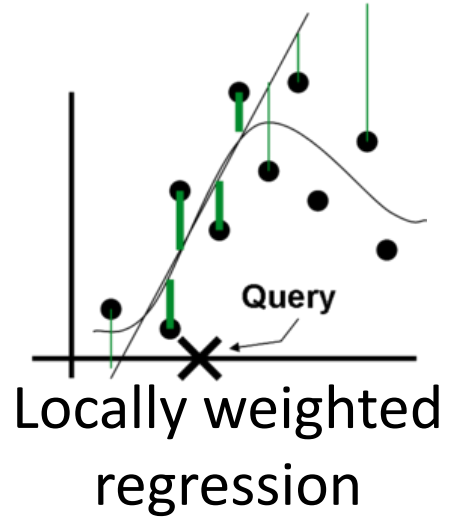
- A measure of how unfaithful g is in approximating f in the locality defined by π_x
- Get an explanation model $\xi(x)$ by optimizing it.

Sampling example - images

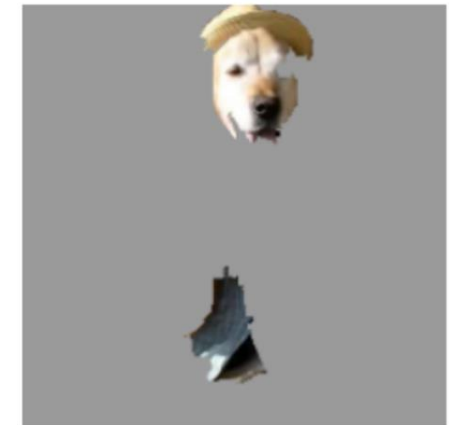


Perturbed Instances	P(Labrador)
	<div data-bbox="1200 511 1421 594"></div> 0.92
	<div data-bbox="1200 818 1225 901"></div> 0.001
	<div data-bbox="1220 1136 1289 1219"></div> 0.34

Original Image
 $P(\text{labrador}) = 0.21$



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



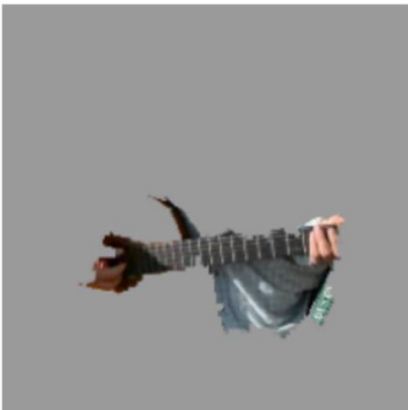
Explanation

Gaining insights from explanations

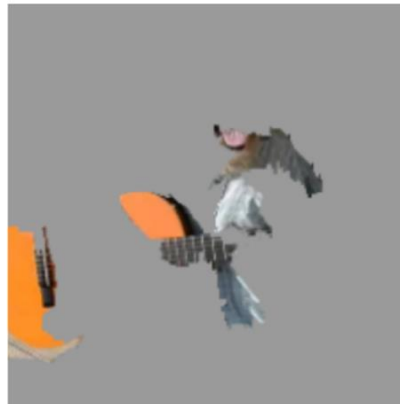
Explaining Google's Inception NN



$$P(\text{🎸}) = 0.32$$



$$P(\text{🎸}) = 0.24$$



$$P(\text{🐶}) = 0.21$$



Train a neural network to predict **wolf** v. **husky**

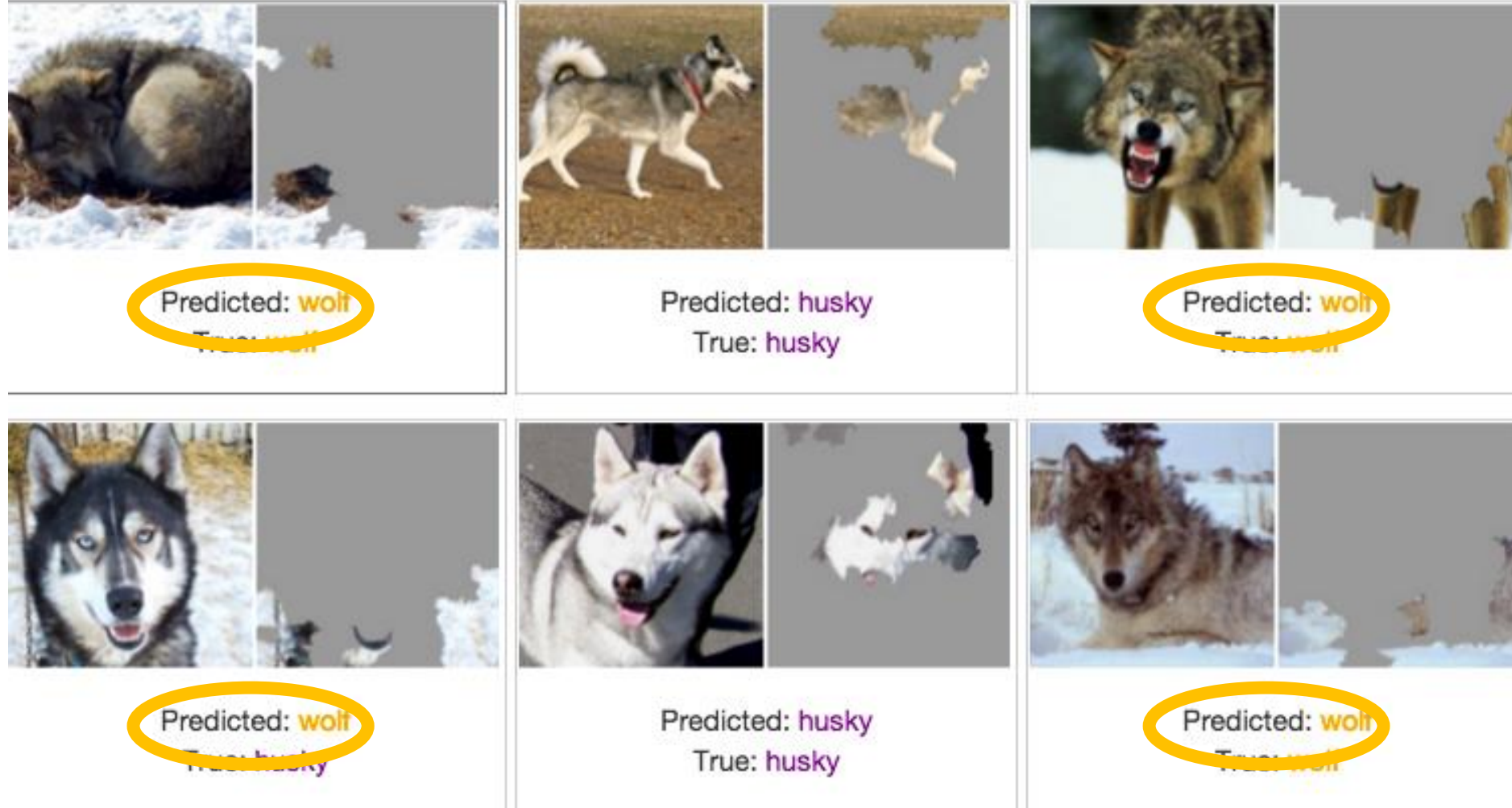


Only 1 mistake!!!

Do you trust this model?

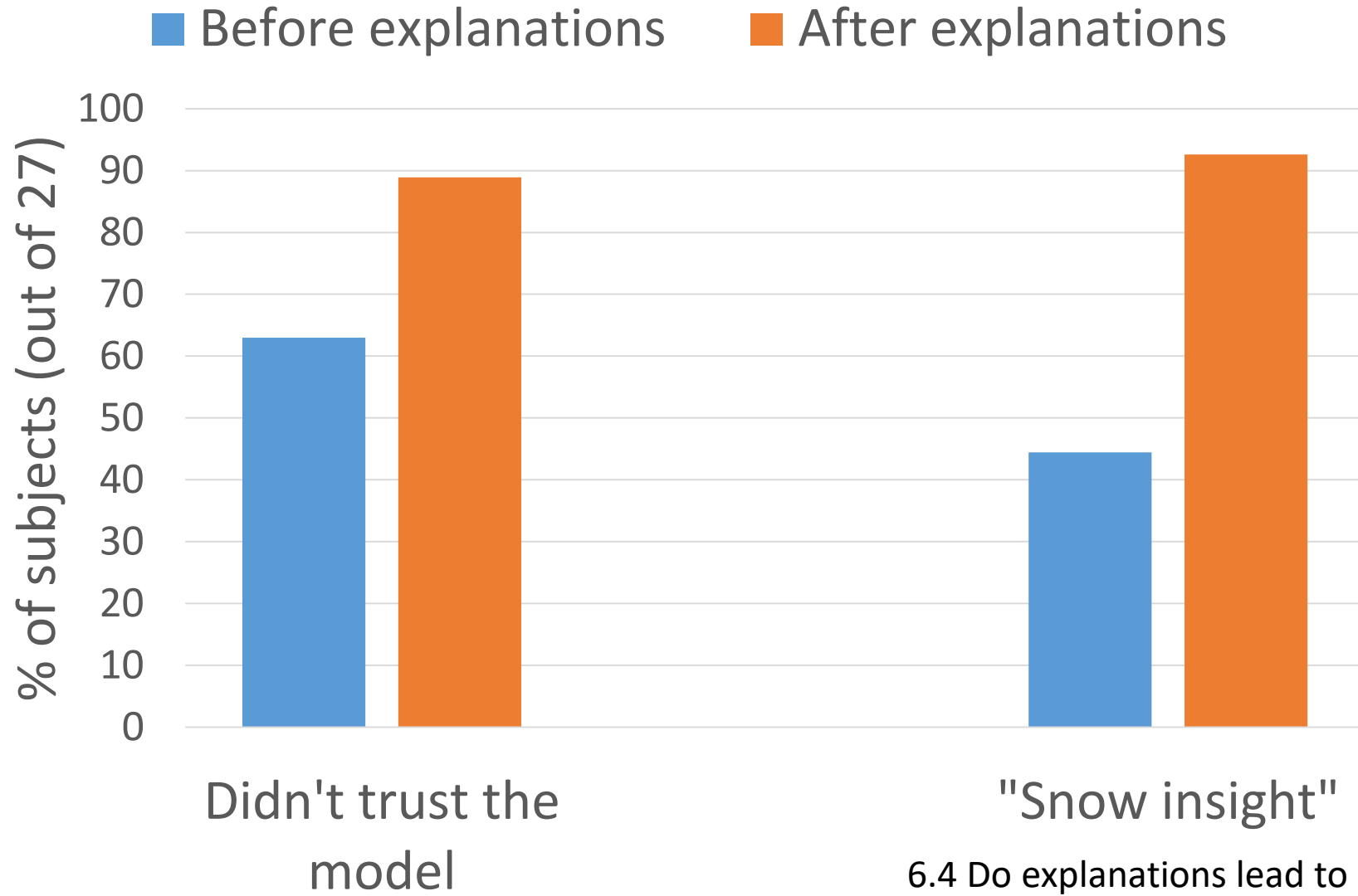
How does it distinguish between huskies and wolves?

Explanations for neural network prediction



We've built a great snow detector... ☹️

Did machine learning people notice it?



Beyond explaining predictions:

Explaining whole models

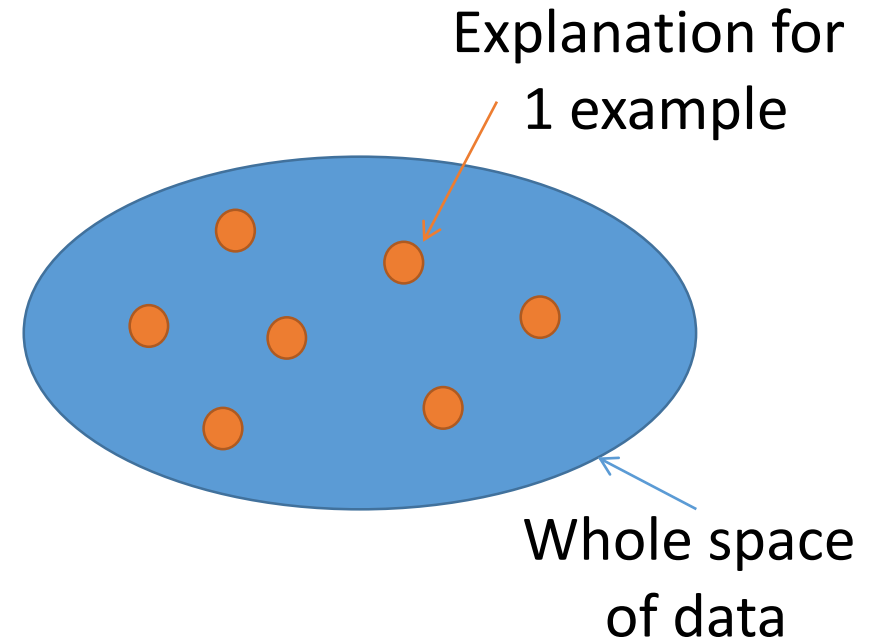
Explaining model behavior for whole space

Explaining 1 prediction describes local behavior of model

Pick k predictions to describe overall behavior

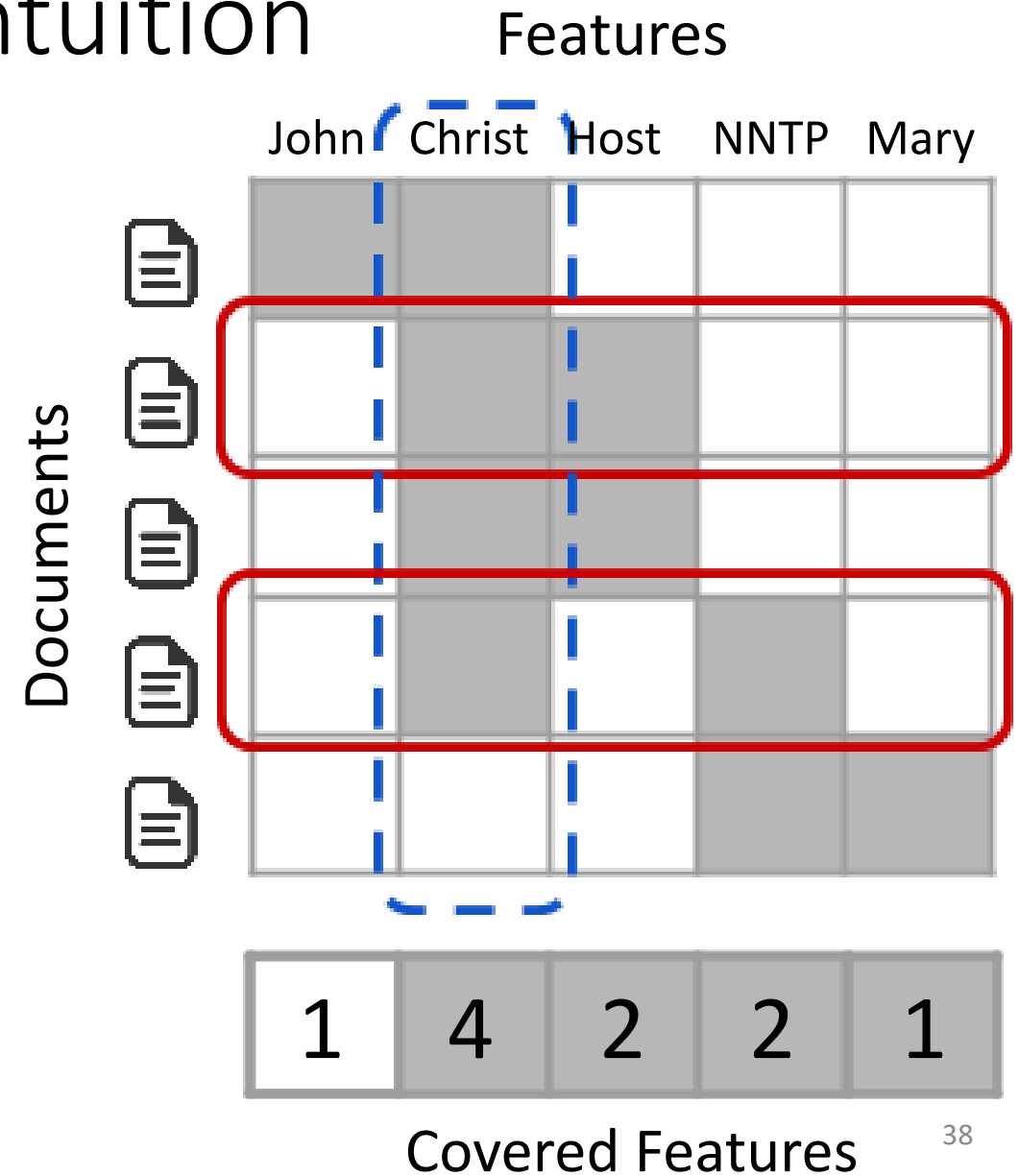
Indispensable to avoid redundancy in explanations

Submodular function optimization provides diverse set of explanations



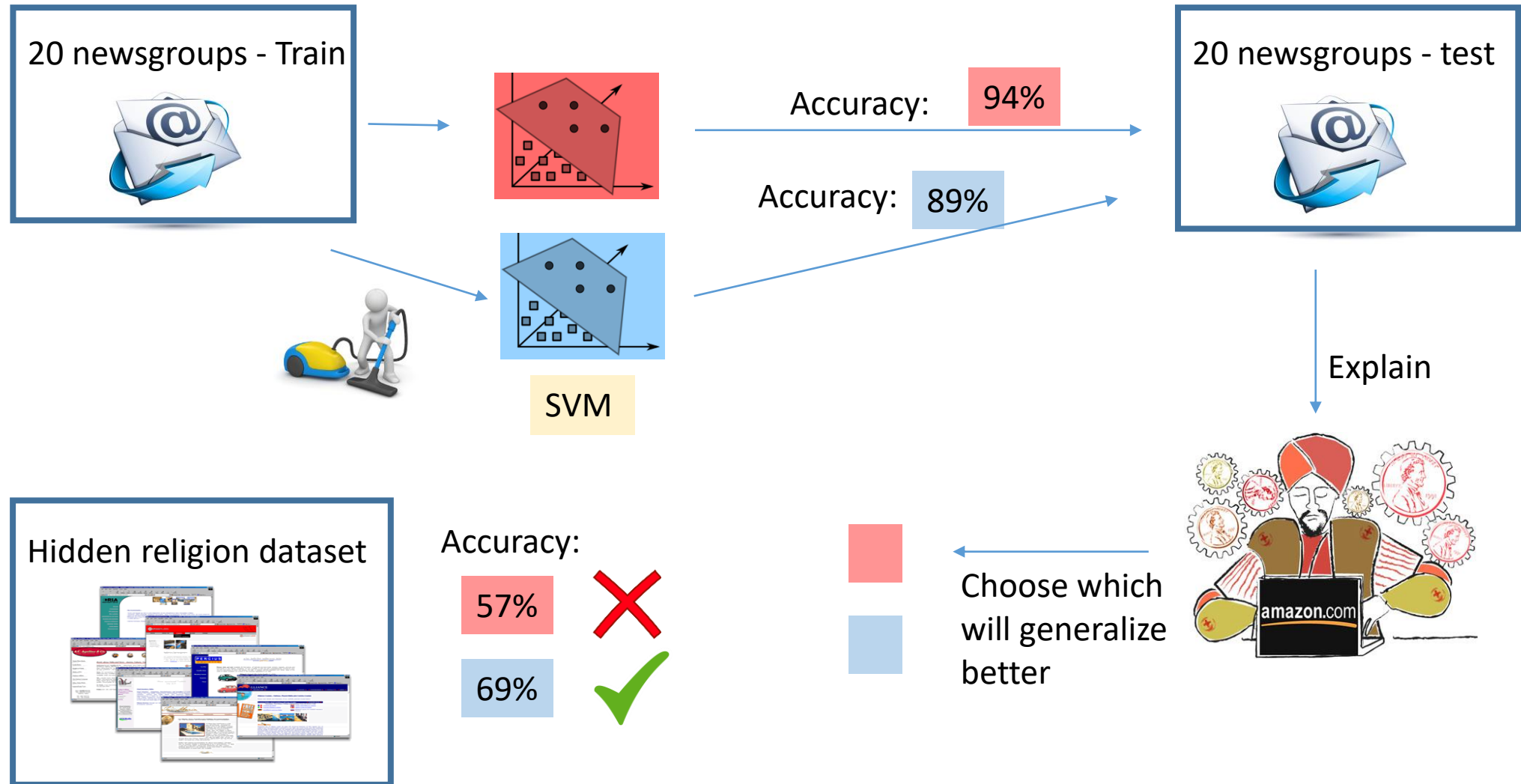
Submodular selection- intuition

1. Representative explanations
2. Avoid redundancy



Evaluating whole-model explanations


Choosing between competing models





Choosing between competing models

- Ask people on Mechanical Turk which model generalizes better

Example #3 of 6

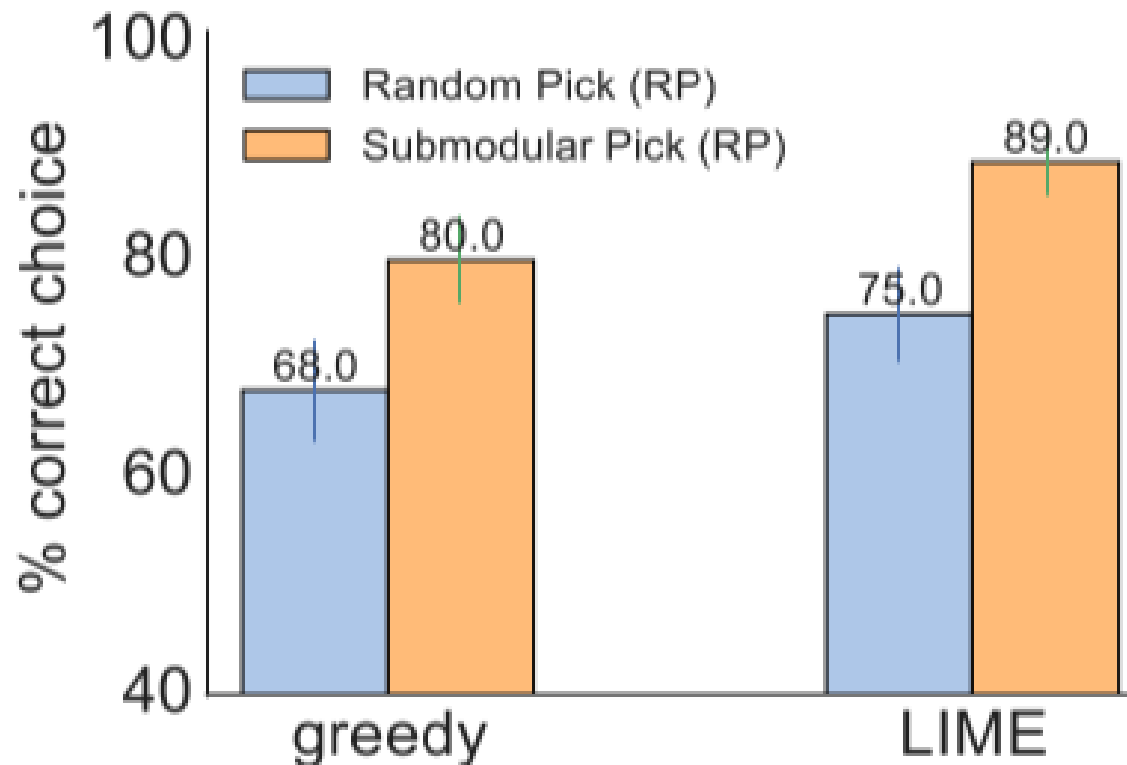
True Class:  Atheism

[Instructions](#) [Previous](#) [Next](#)

Algorithm 1	Algorithm 2																												
<p>Words that A1 considers important:</p> <table border="1"><thead><tr><th>Word</th><th>Importance</th></tr></thead><tbody><tr><td>GOD</td><td>High</td></tr><tr><td>mean</td><td>Medium</td></tr><tr><td>anyone</td><td>Low</td></tr><tr><td>this</td><td>Low</td></tr><tr><td>Koresh</td><td>Medium</td></tr><tr><td>through</td><td>Low</td></tr></tbody></table>	Word	Importance	GOD	High	mean	Medium	anyone	Low	this	Low	Koresh	Medium	through	Low	<p>Words that A2 considers important:</p> <table border="1"><thead><tr><th>Word</th><th>Importance</th></tr></thead><tbody><tr><td>Posting</td><td>High</td></tr><tr><td>Host</td><td>High</td></tr><tr><td>Re</td><td>Medium</td></tr><tr><td>by</td><td>Low</td></tr><tr><td>in</td><td>Low</td></tr><tr><td>Nntp</td><td>Medium</td></tr></tbody></table>	Word	Importance	Posting	High	Host	High	Re	Medium	by	Low	in	Low	Nntp	Medium
Word	Importance																												
GOD	High																												
mean	Medium																												
anyone	Low																												
this	Low																												
Koresh	Medium																												
through	Low																												
Word	Importance																												
Posting	High																												
Host	High																												
Re	Medium																												
by	Low																												
in	Low																												
Nntp	Medium																												
<p>Predicted:</p> <p> Atheism</p> <p>Prediction correct:</p> <p>✓</p>	<p>Predicted:</p> <p> Atheism</p> <p>Prediction correct:</p> <p>✓</p>																												
<p>Document</p> <p>From: pauld@verdix.com (Paul Durbin) Subject: Re: DAVID CORESH IS! GOD! Nntp-Posting-Host: sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>	<p>Document</p> <p>From: pauld@verdix.com (Paul Durbin) Subject: Re: DAVID CORESH IS! GOD! Nntp-Posting-Host: sarge.hq.verdix.com Organization: Verdix Corp Lines: 8</p>																												

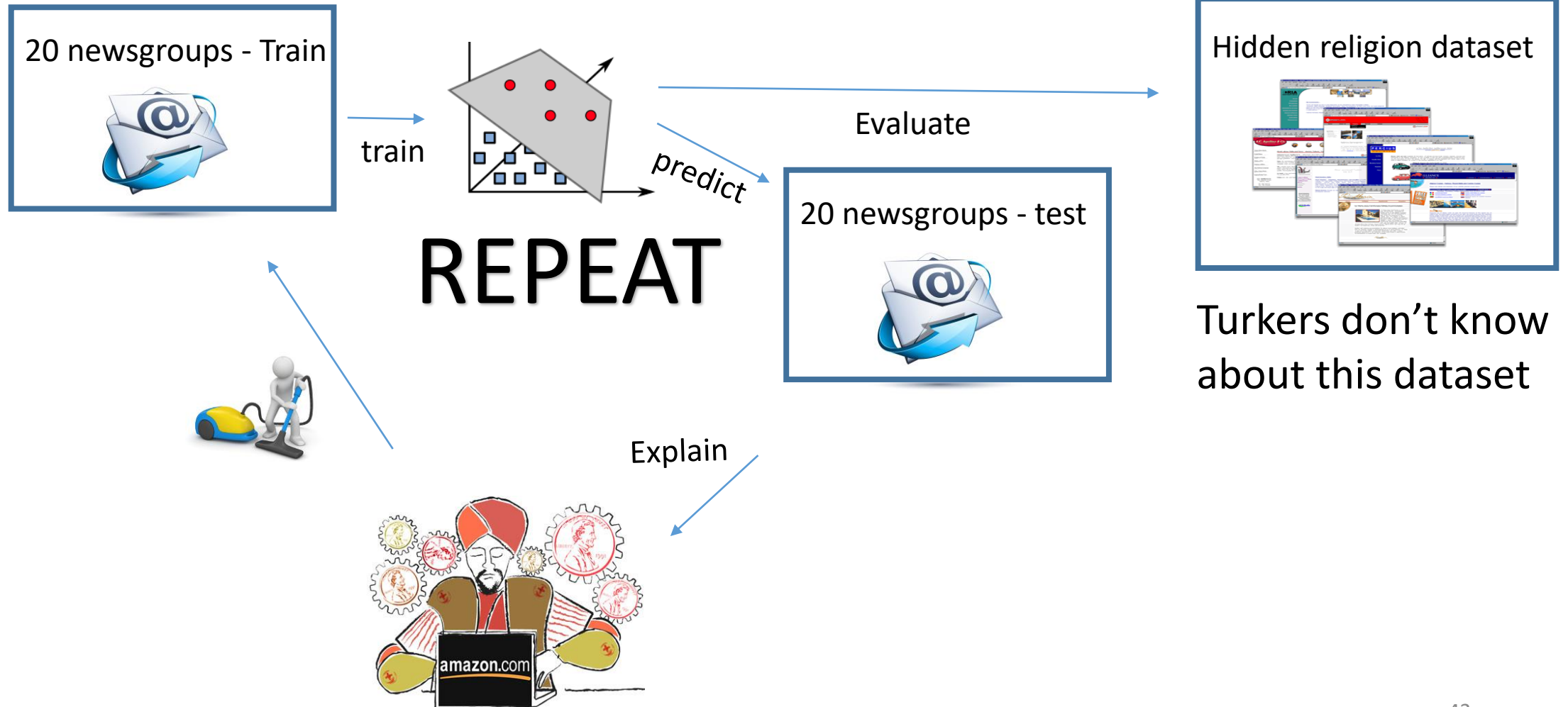
Choosing between competing models

89% of Mechanical Turkers
identify more trustworthy model



If we picked based
on accuracy, we
would get it wrong.

Fixing bad classifiers



Fixing bad classifiers

- Turkers click on 'useless' words for the task in each round

Example #5 of 10

True Class: ● Atheism

[Instructions](#) [Previous](#) [Next](#)

Words that the algorithm considers important.

Host	<div></div>
Posting	<div></div>
NNTP	<div></div>
to	<div></div>
New	<div></div>
Thanks	<div></div>
anyone	<div></div>
email	<div></div>
not	<div></div>
has	<div></div>

Bar length indicates importance, and color indicates to which topic: Christianity (green) or Atheism (Pink).

Please click on the words (right next to the bars) that you think the algorithm is using incorrectly, because they are not important to distinguish between Atheism and Christianity. They should be red and crossed off after you click them.

Document

From: johnchad@triton.unm.edu (jchadwic)
Subject: Another request for Darwin Fish
Organization: University of New Mexico, Albuquerque
Lines: 11
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

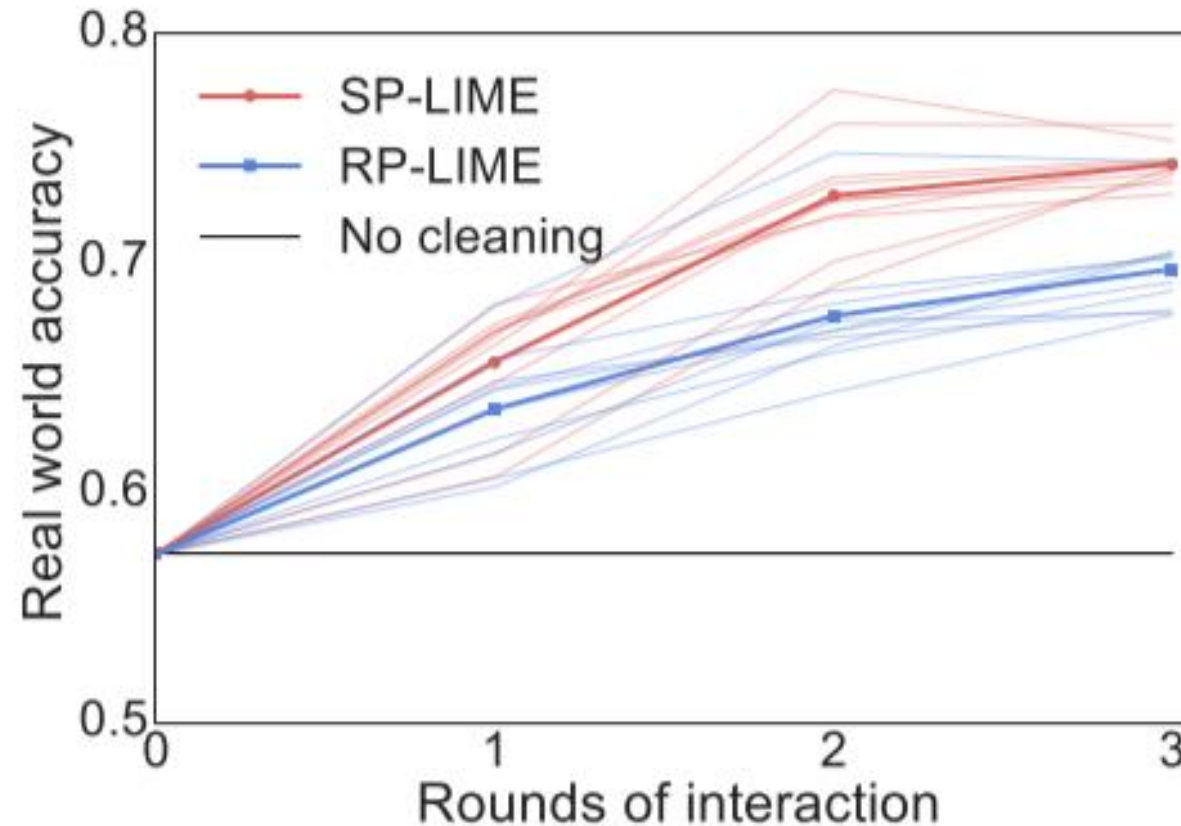
There have been some notes recently asking where to obtain the DARWIN fish.
This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.

Thanks,

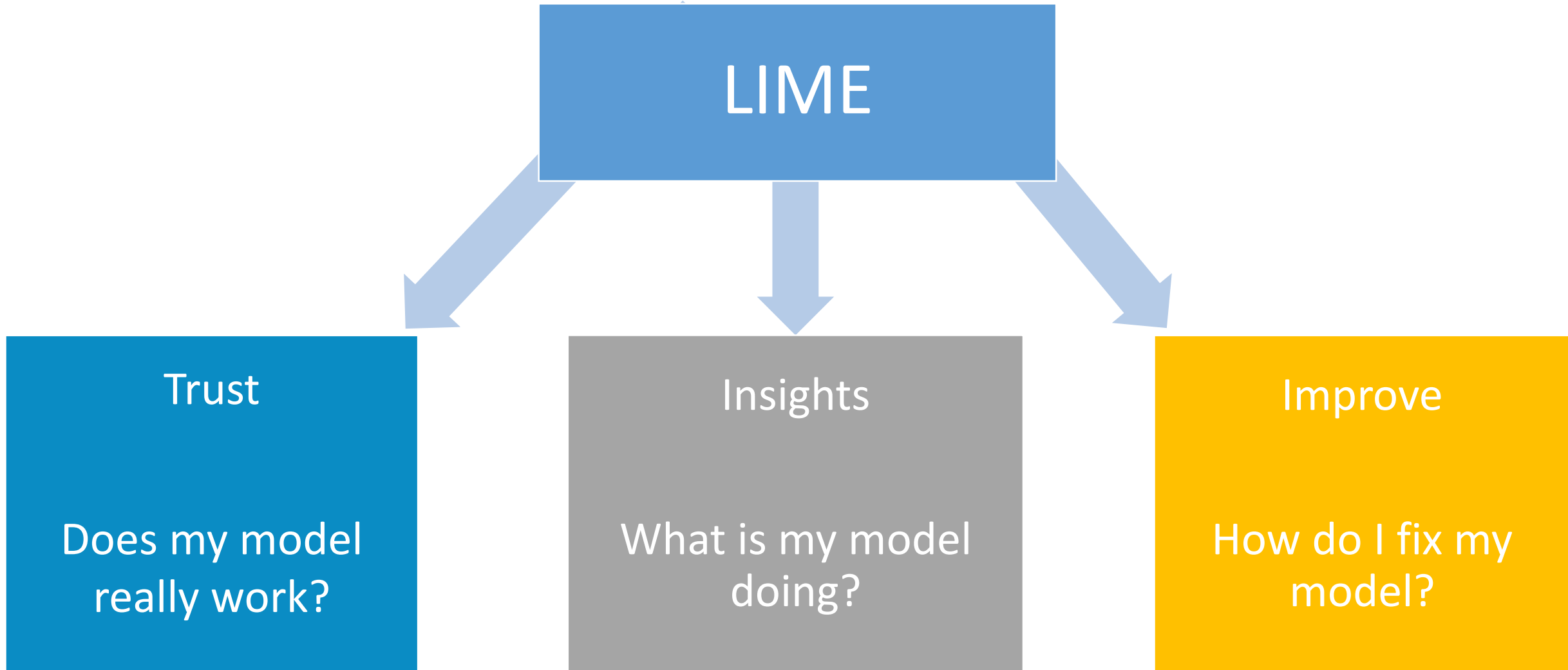
john chadwick
johnchad@triton.unm.edu
or

Fixing bad classifiers

Mechanical Turkers can do
'feature engineering' really well!



Conclusion



Open source project: <https://github.com/marcotcr/lime>