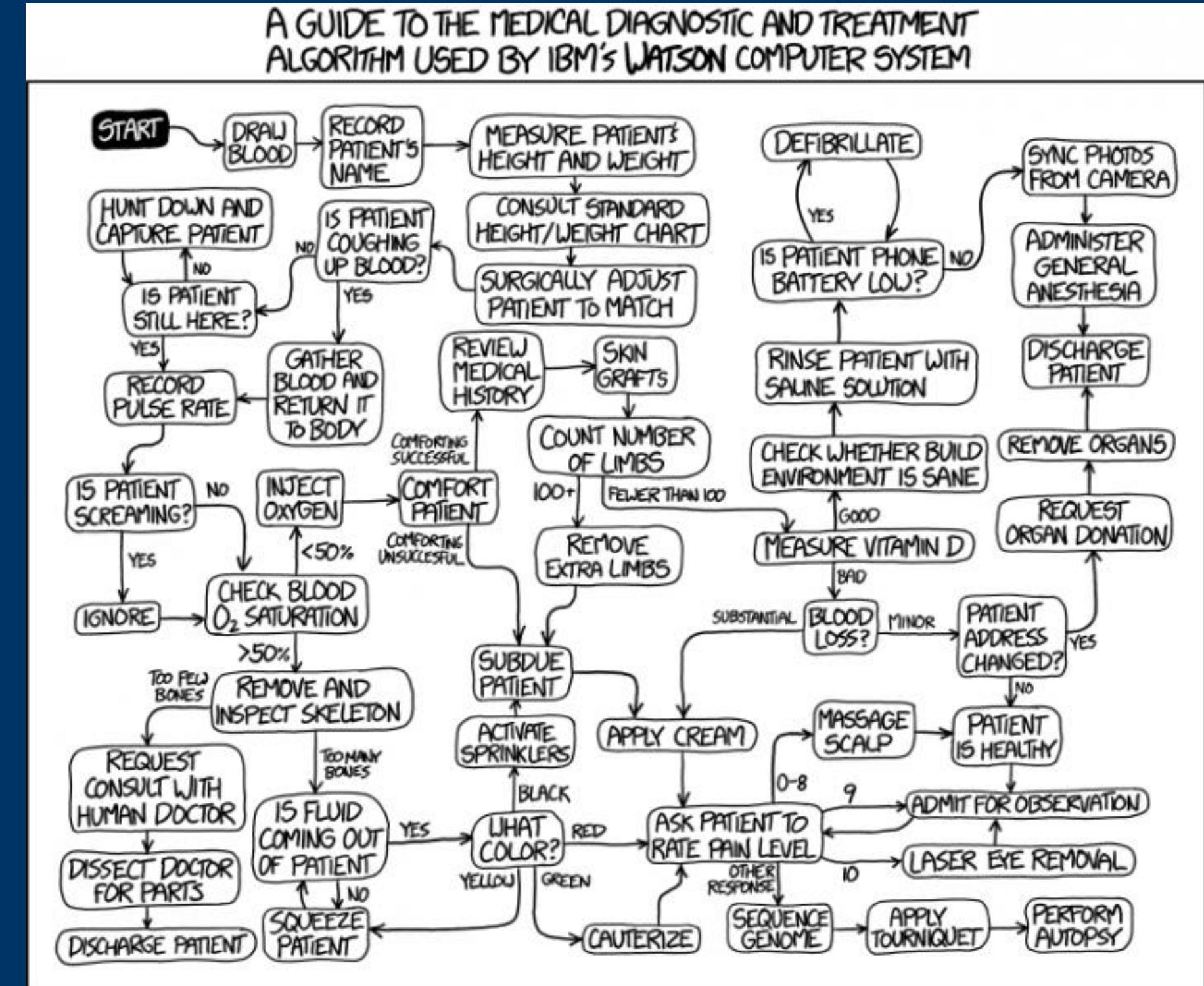


# Excuse Me, Why Has Your AI Denied My Loan?

## A Look at Machine Learning Model Explainability

Sorin Pește  
Cloud Solution Architect, Data & AI @ Microsoft



# CODECAMP\_BUCHAREST

/ 9 November 2019

/ JW Marriott Grand Hotel Bucharest

## Global partners

Cognizant  
Softvision



## Diamond partners

METRO  
SYSTEMS



UBISOFT



tremend



pipedrive



## Platinum partners

stefanini

SYSTEMATIC

MindGeek



## Gold partners

### Entrepreneurship



SPINLAB THE HHL ACCELERATOR

NESCAFÉ

vitamin aqua



.exposé



## Liked By

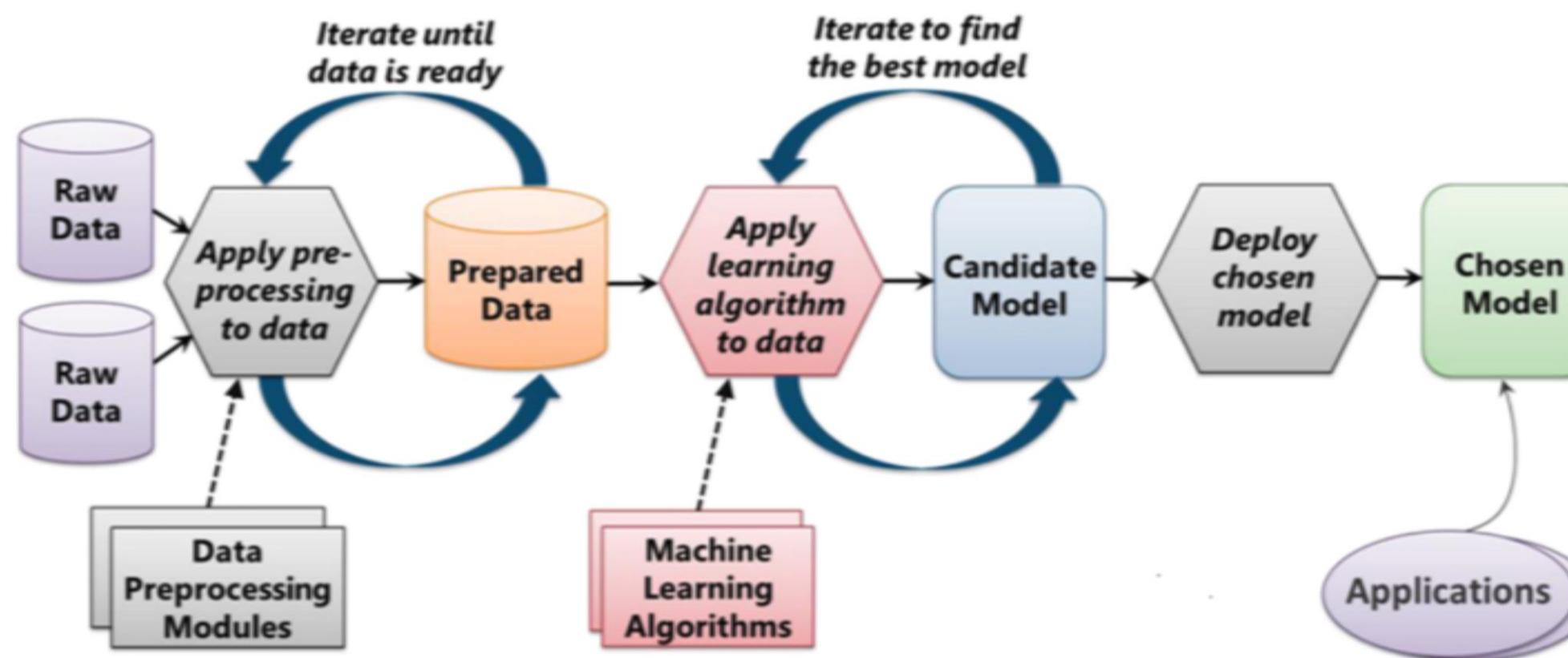


# Demo Code

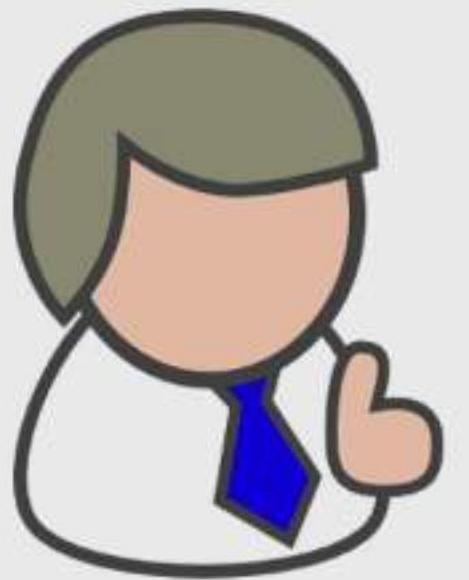
<https://github.com/neaorin/Explaining-ML-Models/>

Why Should You Care About Explainability  
In Machine Learning?

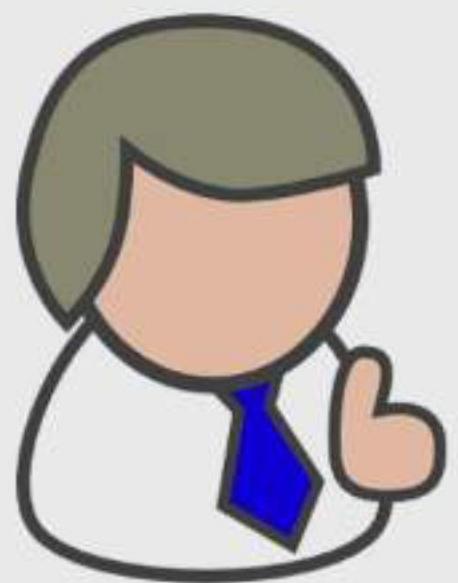
# The Machine Learning Process



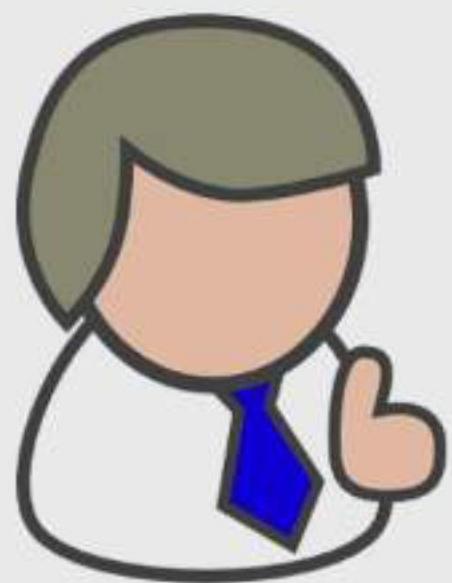
From "Introduction to Microsoft Azure" by David Chappell



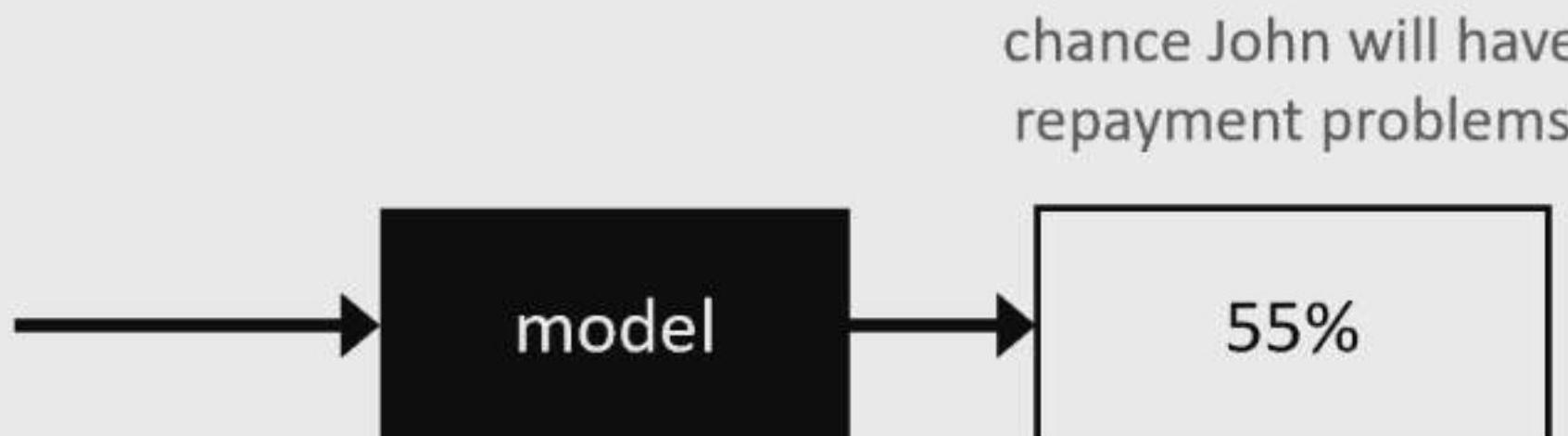
John, a bank customer

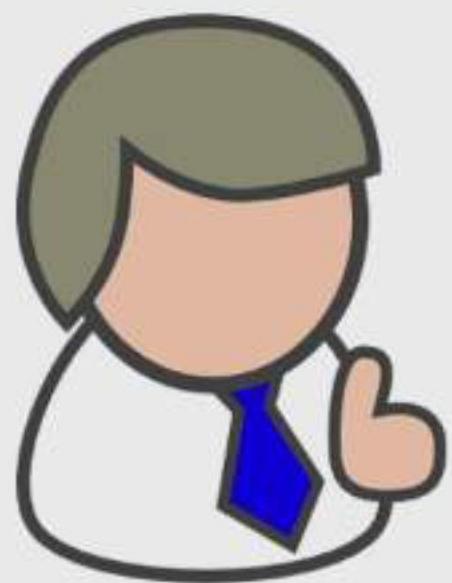


John, a bank customer

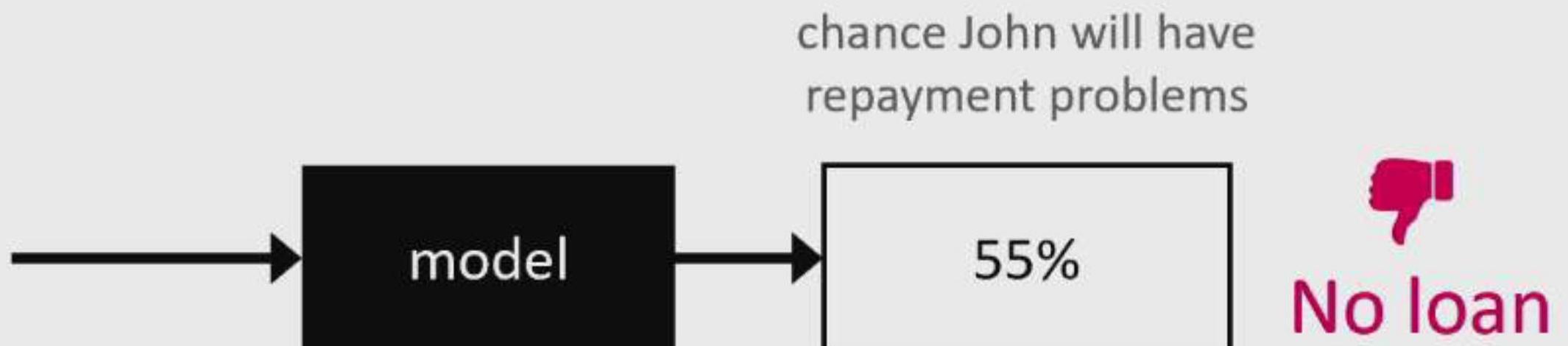


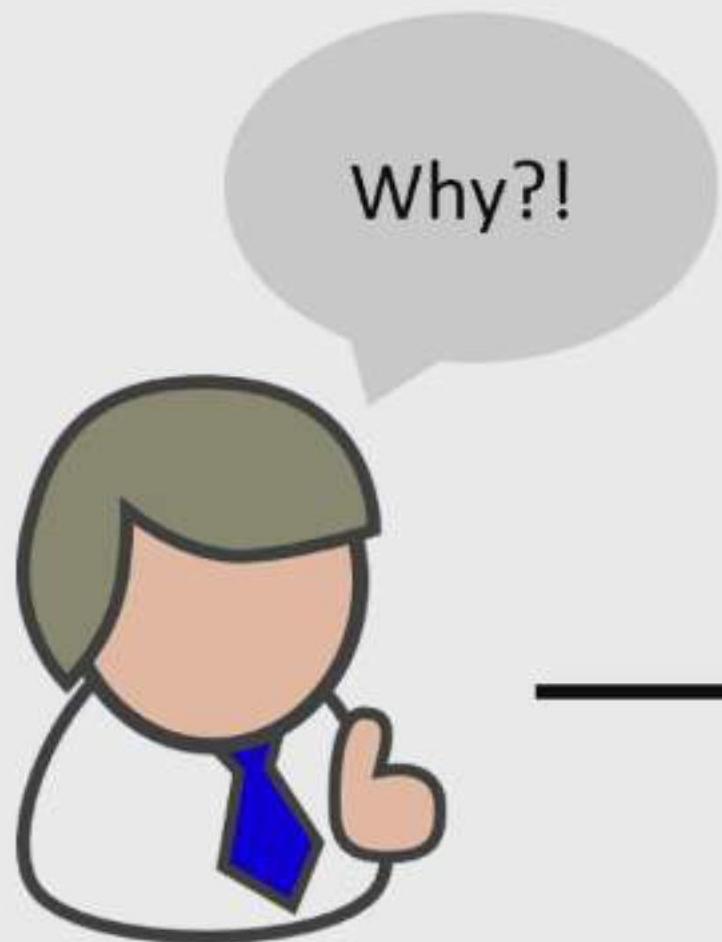
John, a bank customer



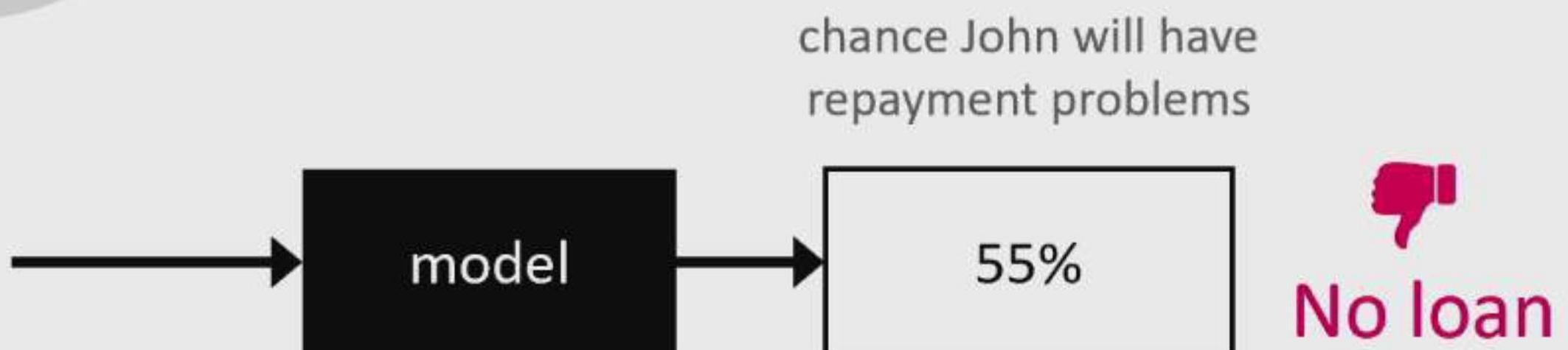


John, a bank customer





John, a bank customer



chance John will have  
repayment problems

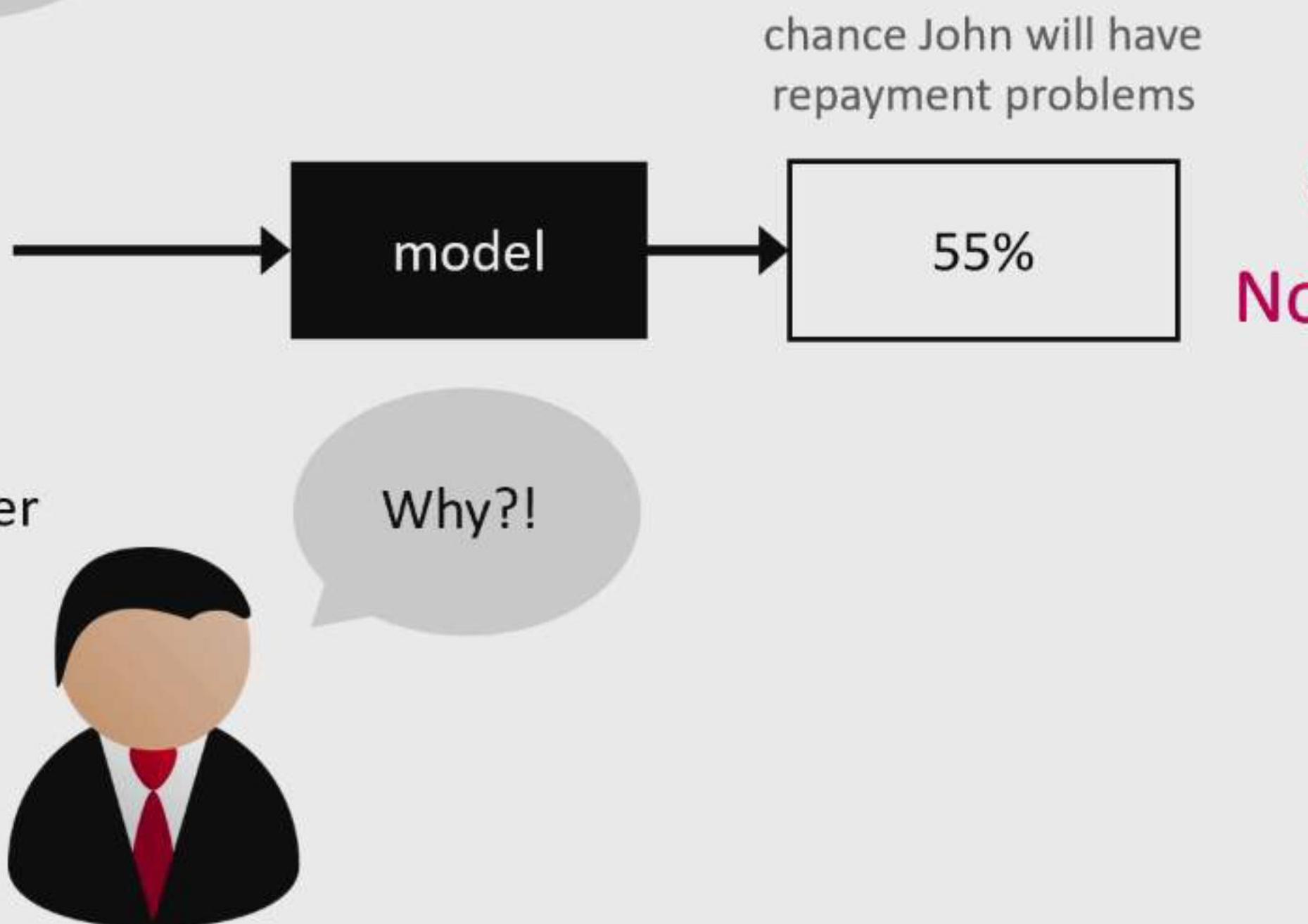
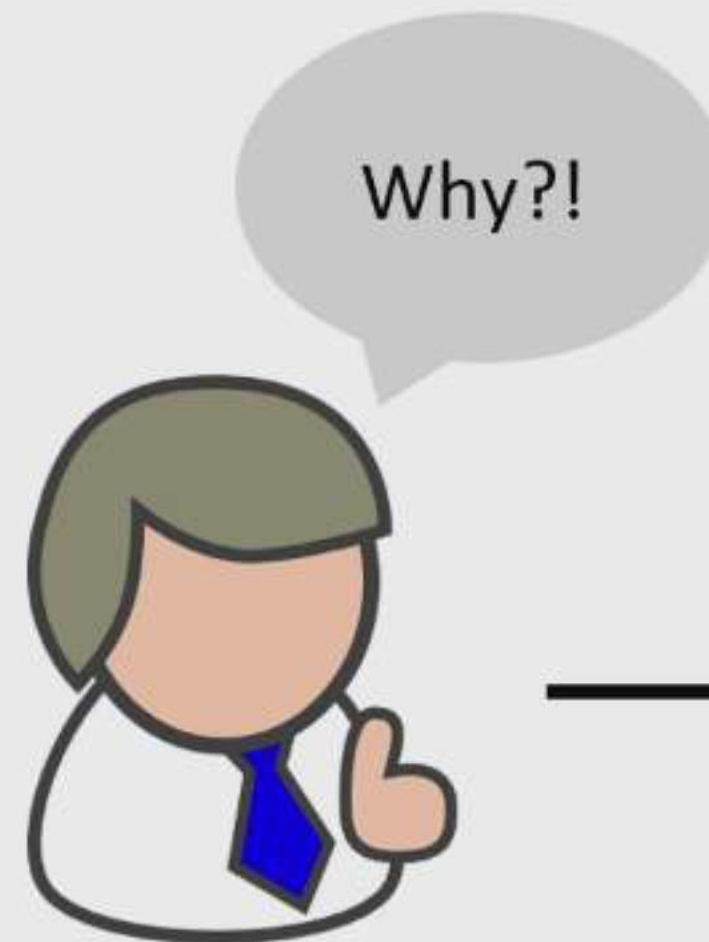
55%



No loan

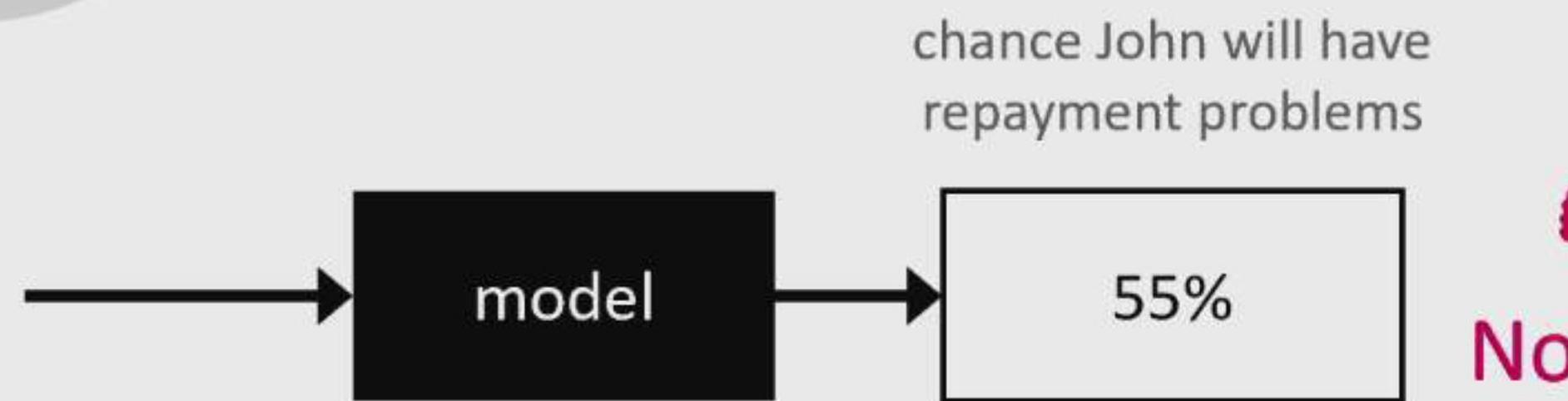
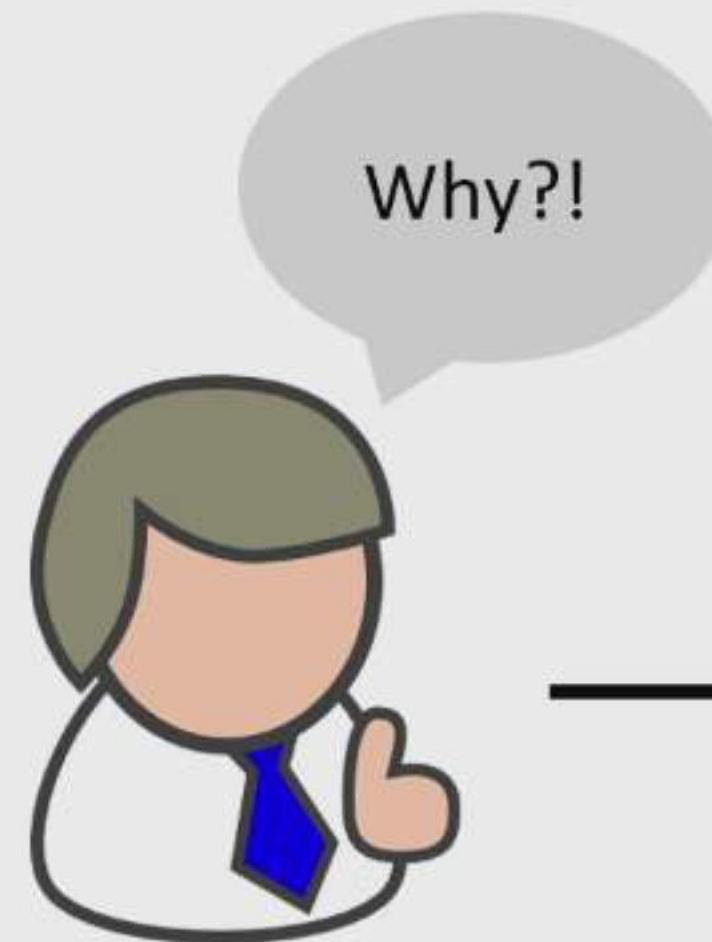


John, a bank customer

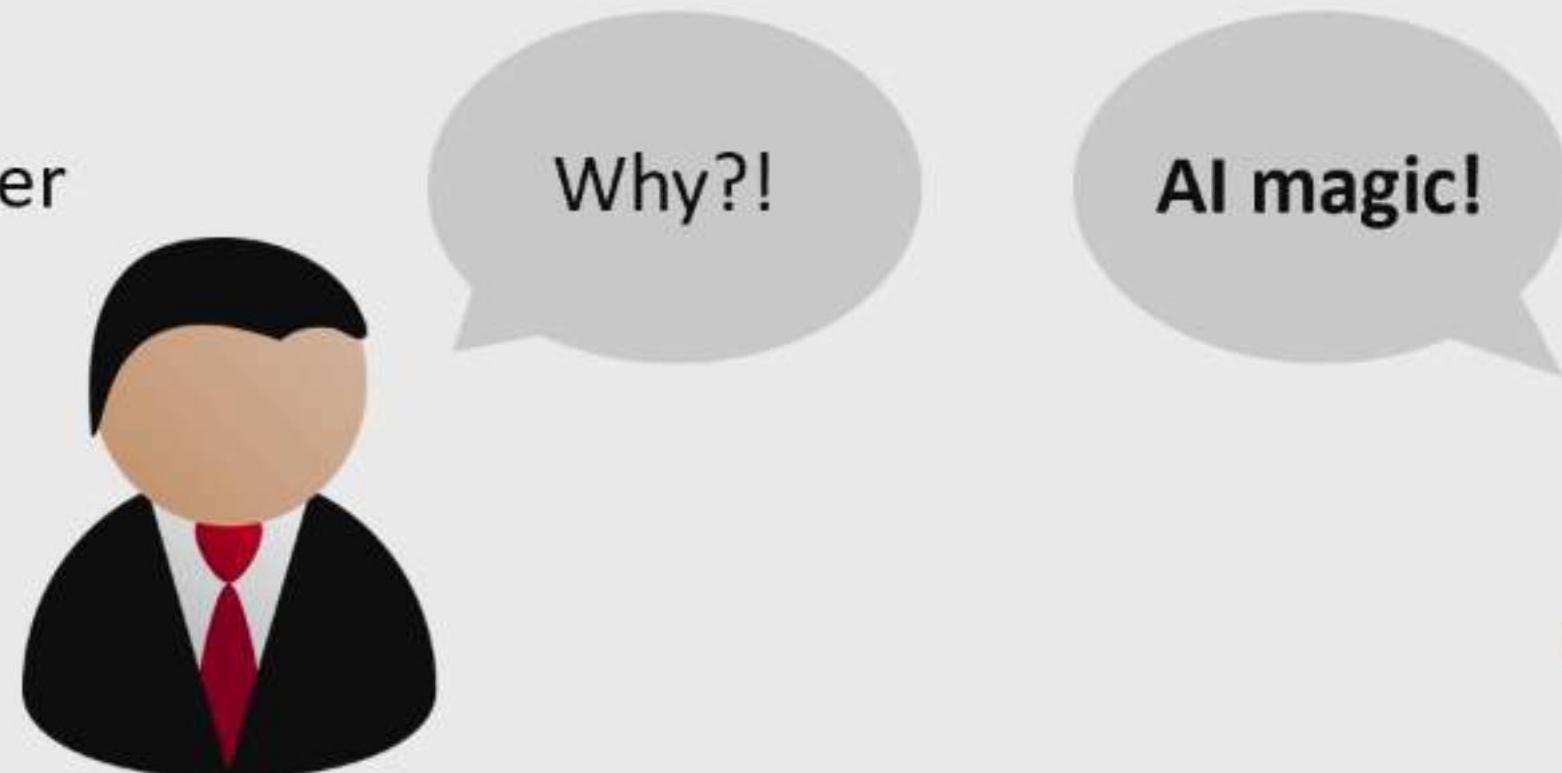




John, a bank customer



No loan



BUSINESS NEWS OCTOBER 10, 2018 / 6:12 AM / A YEAR AGO

# Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin

8 MIN READ



SAN FRANCISCO (Reuters) - Amazon.com Inc's ([AMZN.O](#)) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

Artificial Intelligence

# Military artificial intelligence can be easily and dangerously fooled

AI warfare is beginning to dominate military strategy in the US and China, but is the technology ready?

by Will Knight

Oct 21, 2019

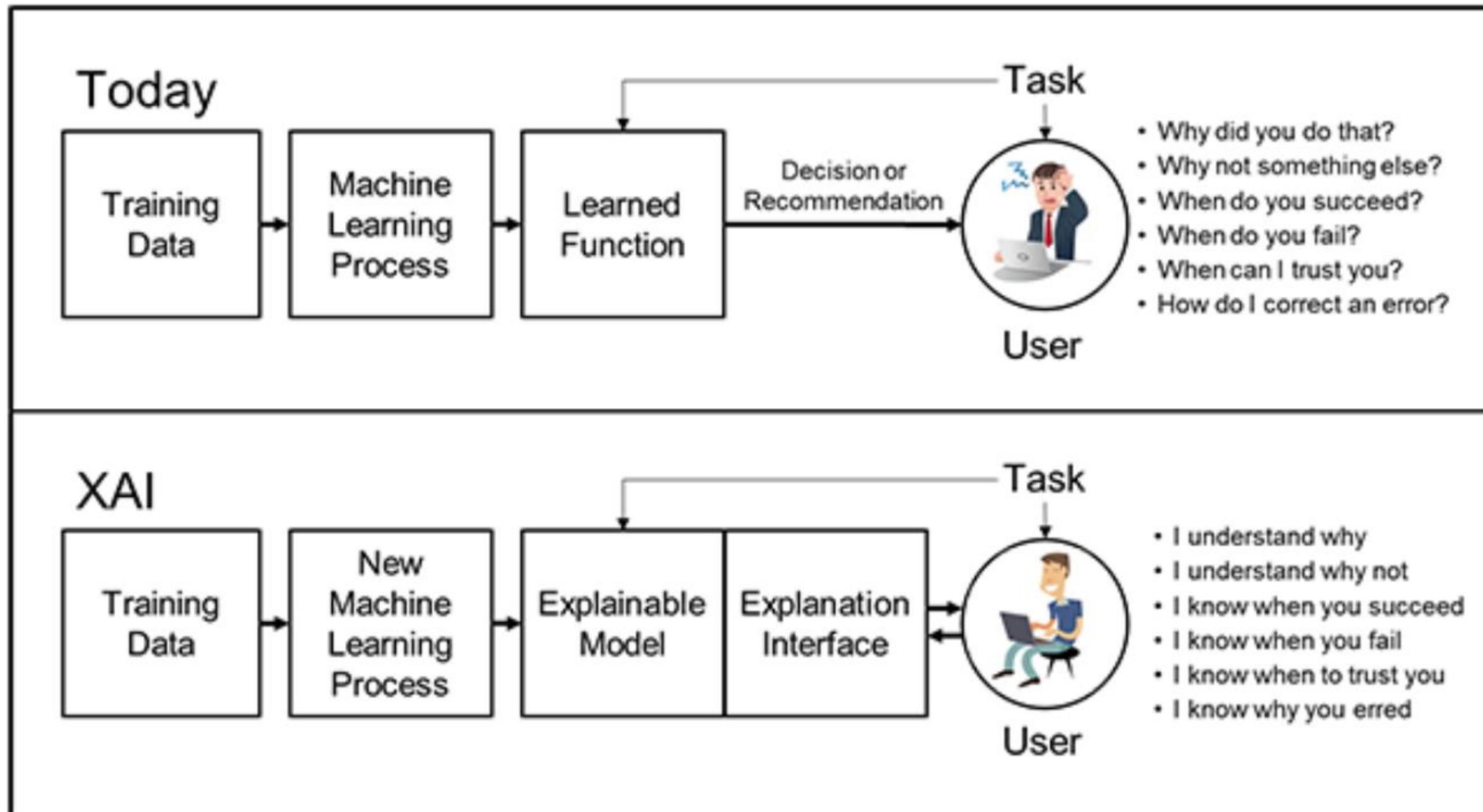
**Last March, Chinese researchers announced an ingenious and potentially devastating attack against one of America's most prized technological assets — a Tesla electric car.**



The team, from the security lab of the Chinese tech giant Tencent, demonstrated several ways to fool the AI algorithms on Tesla's car. By subtly altering the data fed to the car's sensors, the researchers were able to bamboozle and bewilder the artificial intelligence that runs the vehicle.

In one case, a TV screen contained a hidden pattern that

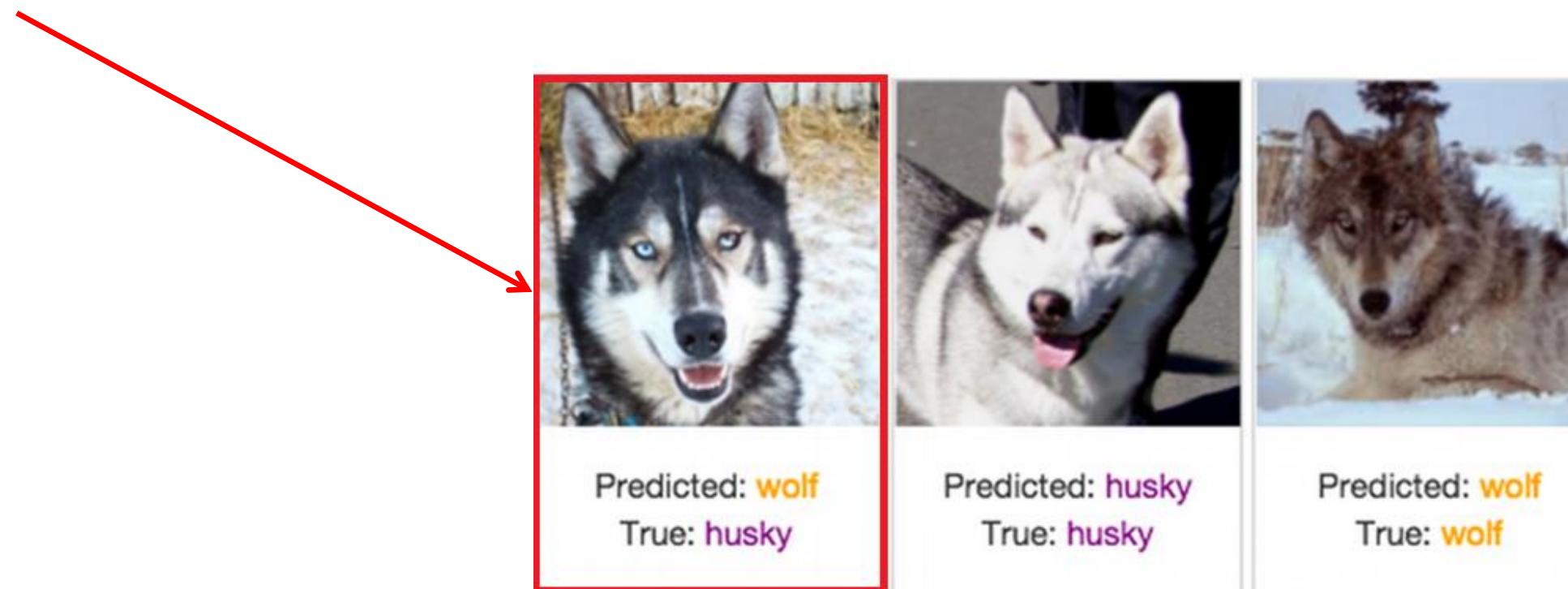
# We Need AI to be Explainable



# Don't Mistake Model Performance for Trustworthiness



Only One Mistake!



# Don't Mistake Model Performance for Trustworthiness



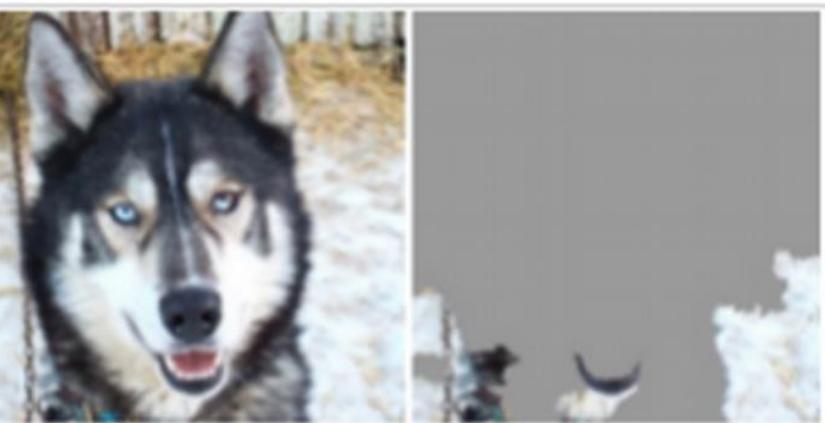
Predicted: **wolf**  
True: **wolf**



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



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



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



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



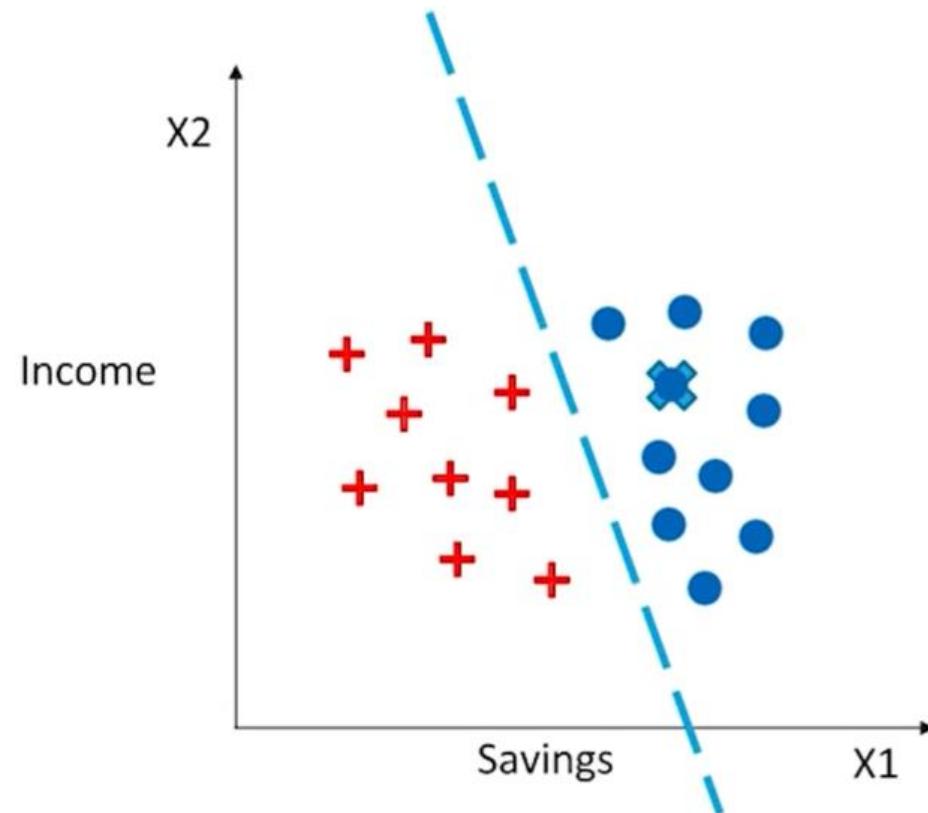
Predicted: **wolf**  
True: **wolf**

	Interpretable	Accurate
Complex model	X	✓
Simple model	✓	X

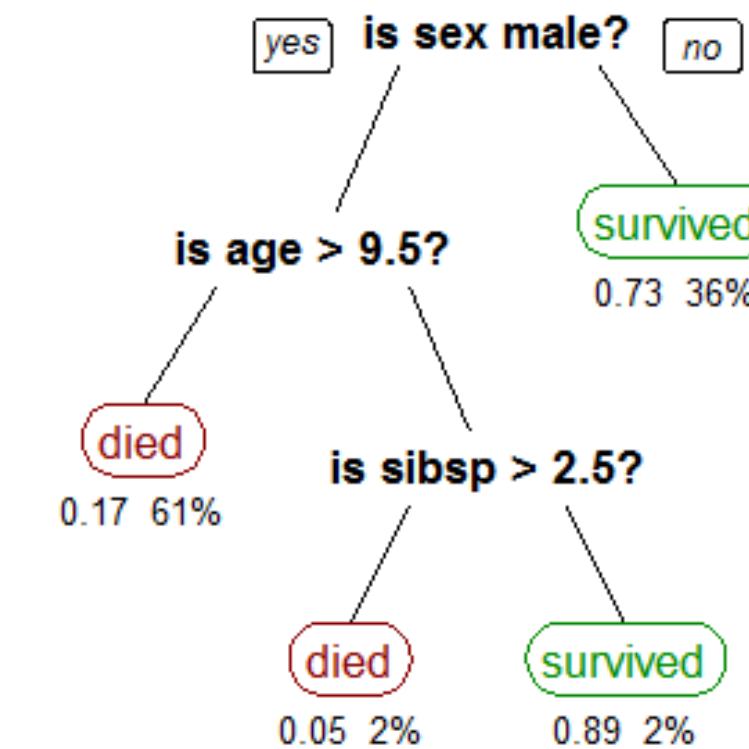
Interpretable or accurate: **choose one.**

# Some Models Are Easy to Interpret

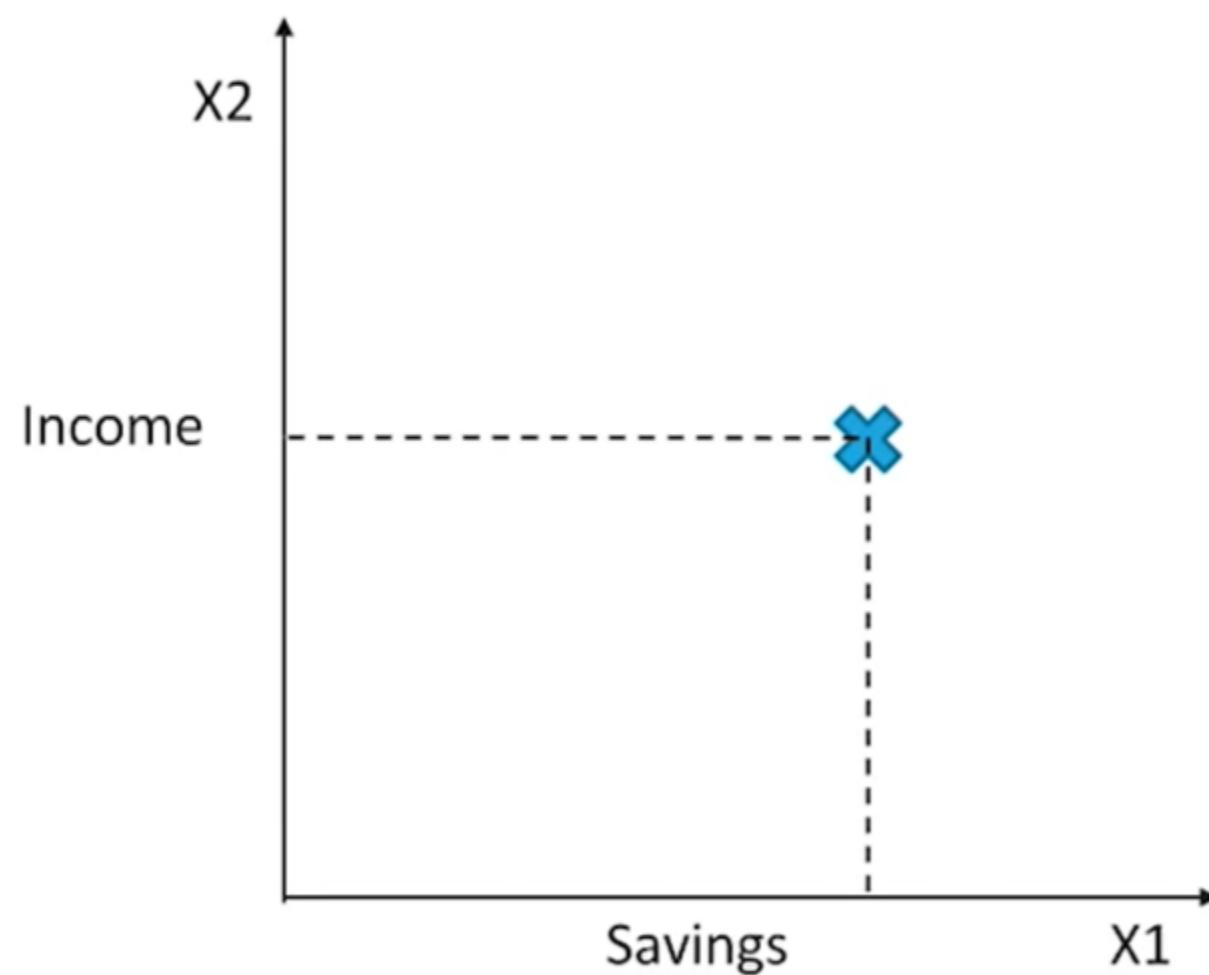
Linear / Logistic Regression



Single Decision Tree

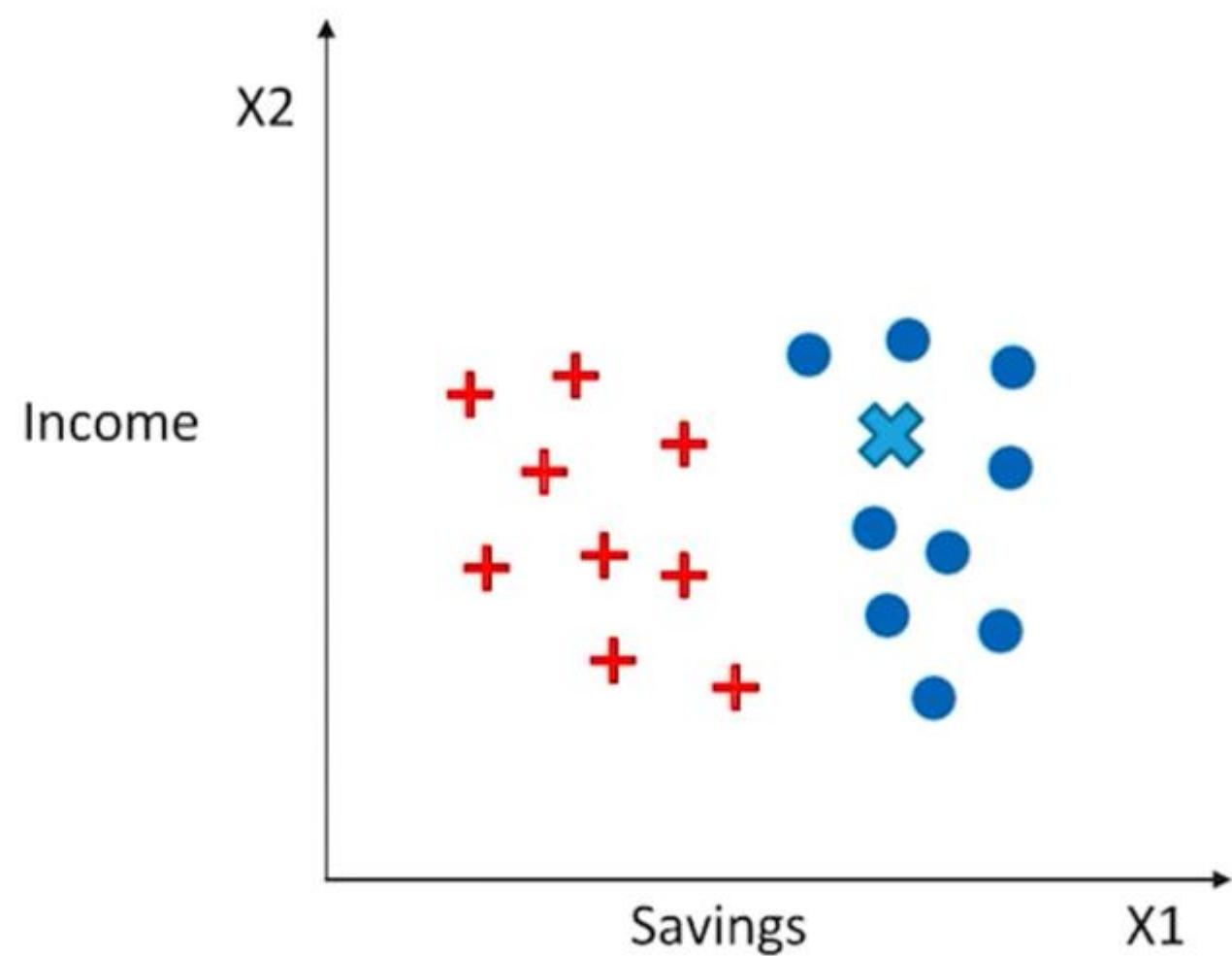


# Will the loan default?



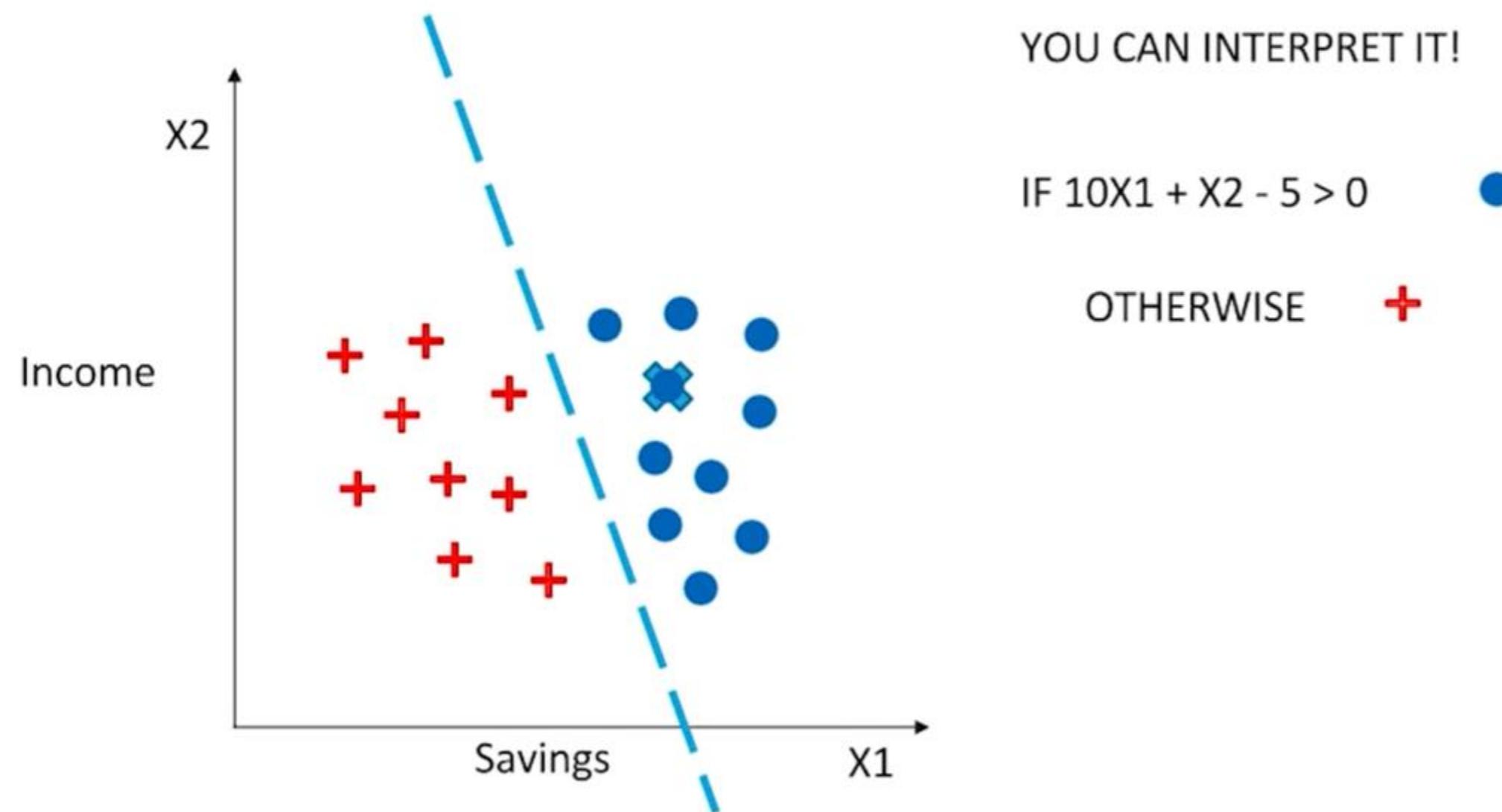
Source: <https://www.youtube.com/watch?v=LAm4QmVaf0E&t=3658s>

# Get Historical Data



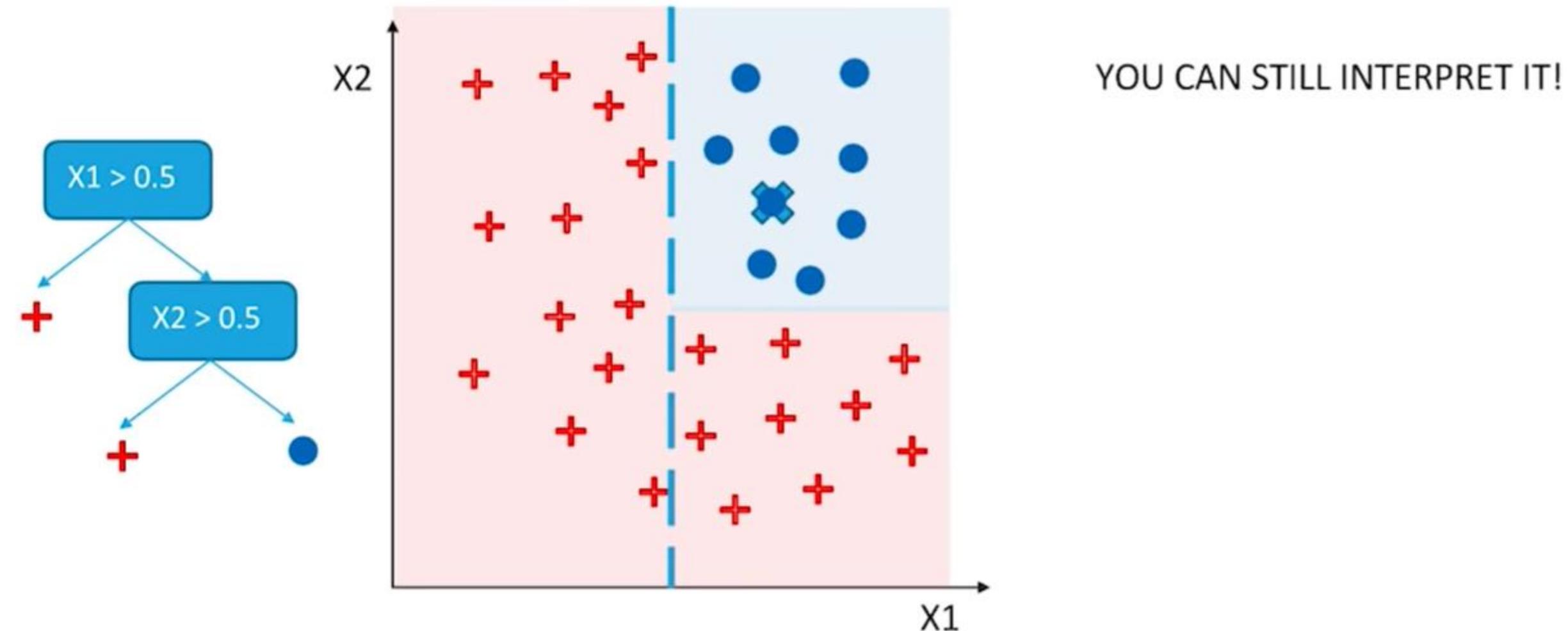
Source: <https://www.youtube.com/watch?v=LAm4QmVaf0E&t=3658s>

# Linear Classifiers



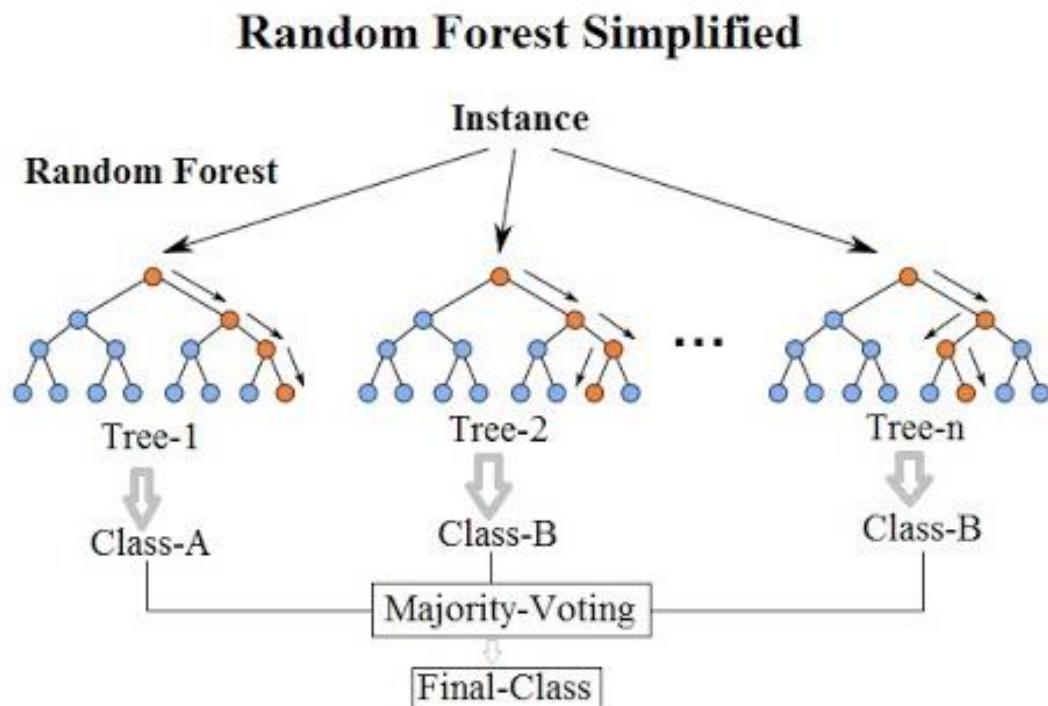
Source: <https://www.youtube.com/watch?v=LAm4QmVaf0E&t=3658s>

# Decision trees

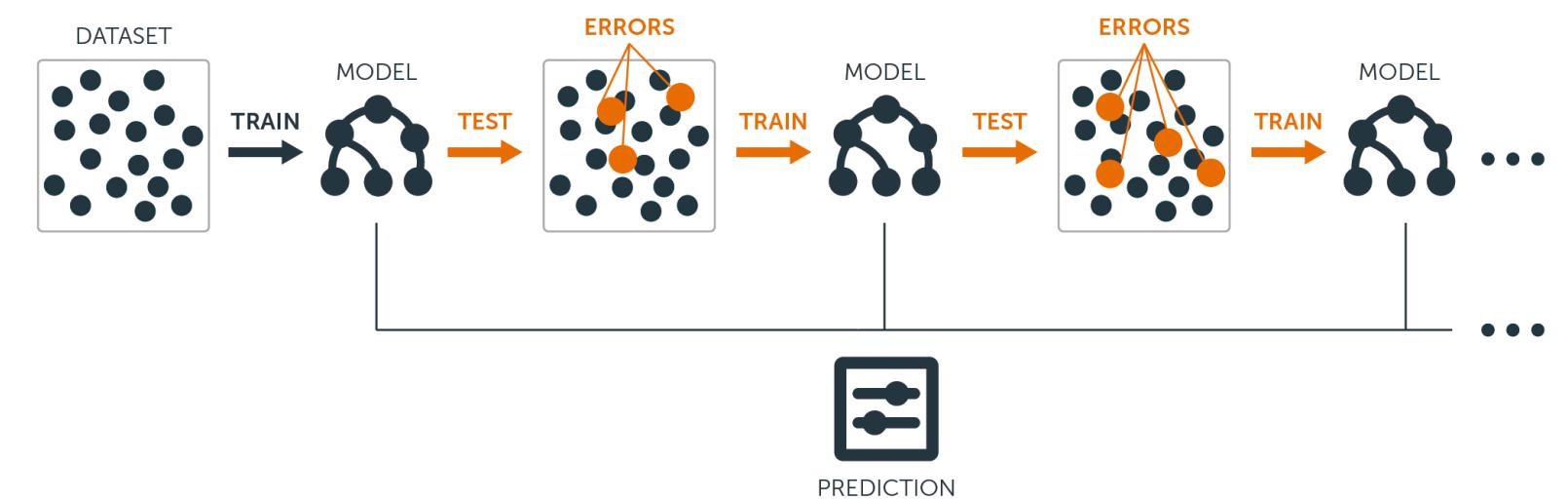


# Some Models Are Harder to Interpret

## Random Forest



## Boosted Trees



# Some Models are REALLY HARD to Interpret

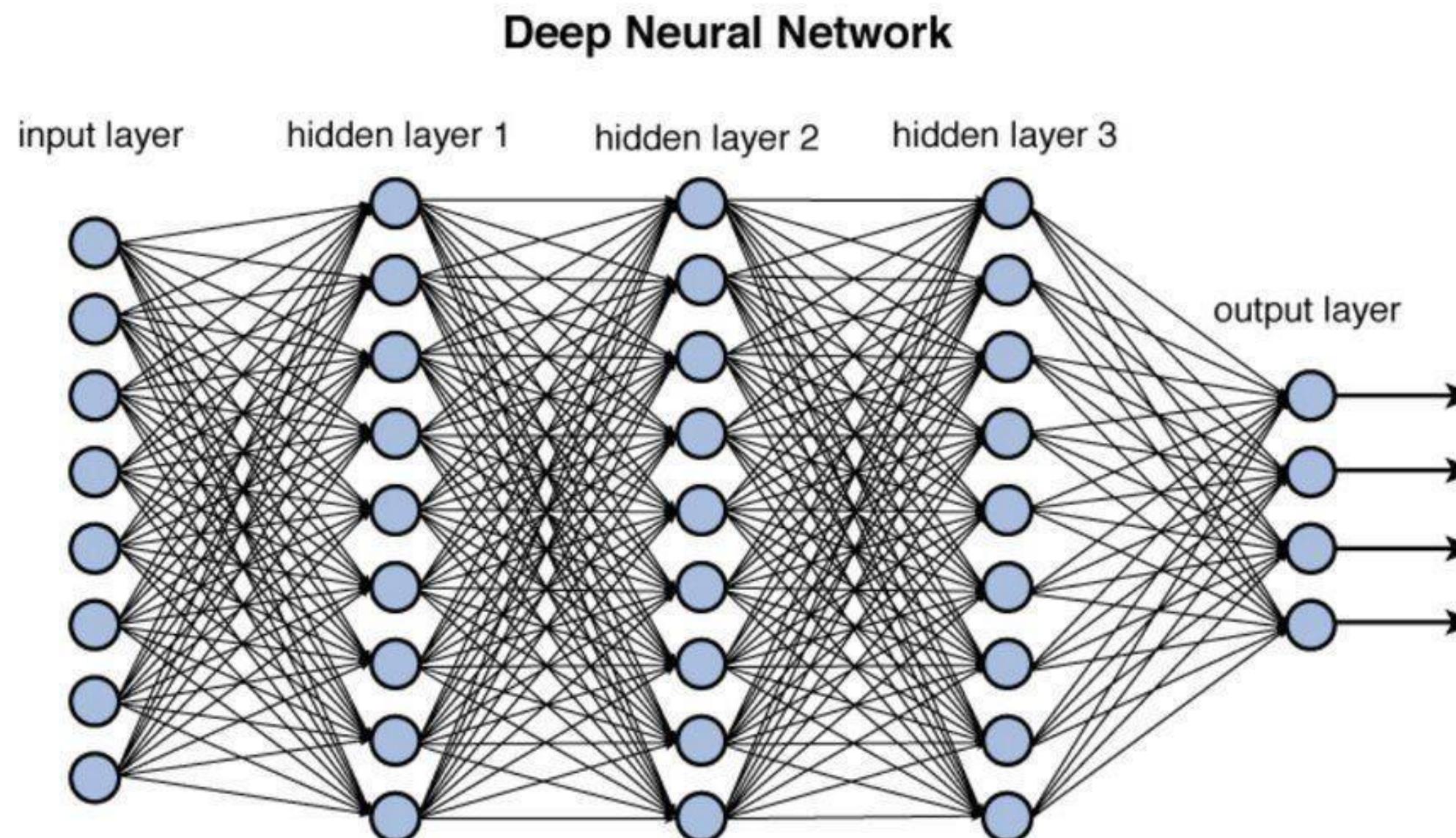


Figure 12.2 Deep network architecture with multiple layers.

# Explainable to Whom?

Data Scientist

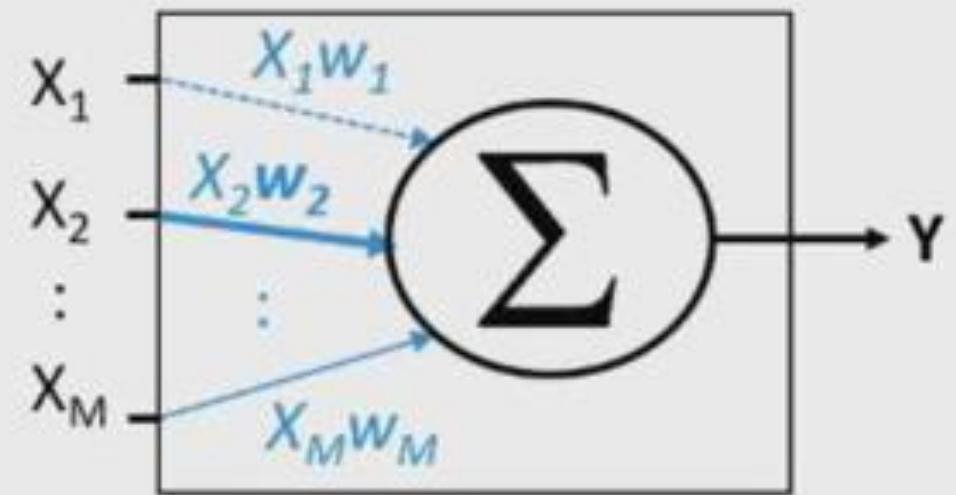


End User



# Types Of Explainers

Glass Box

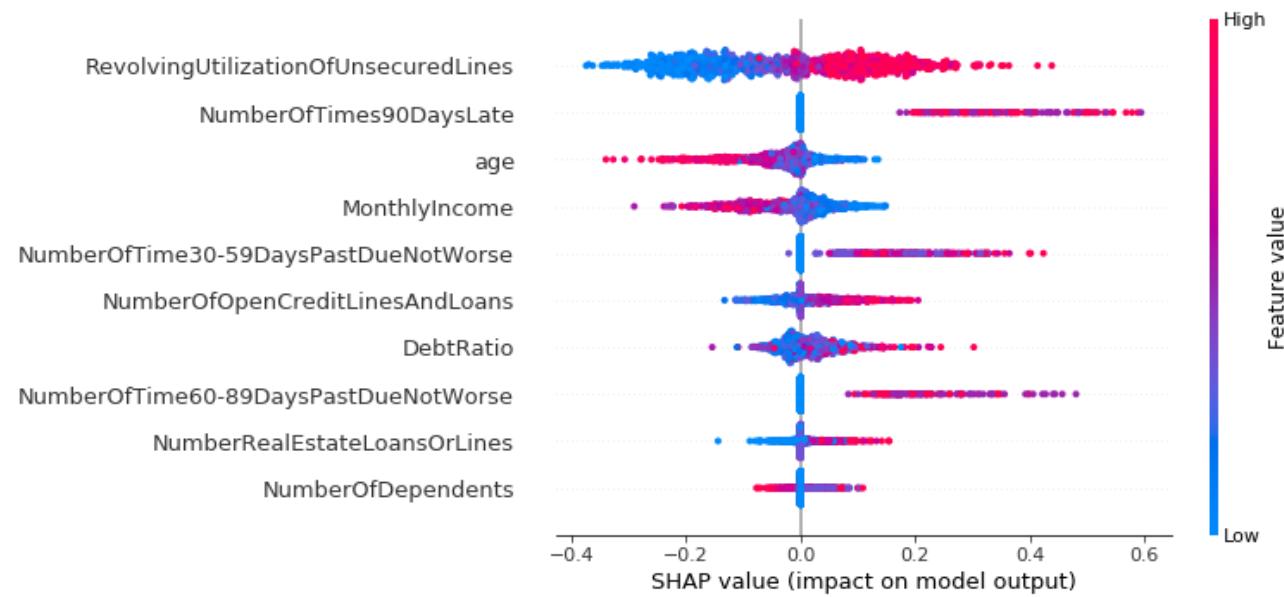


Black Box

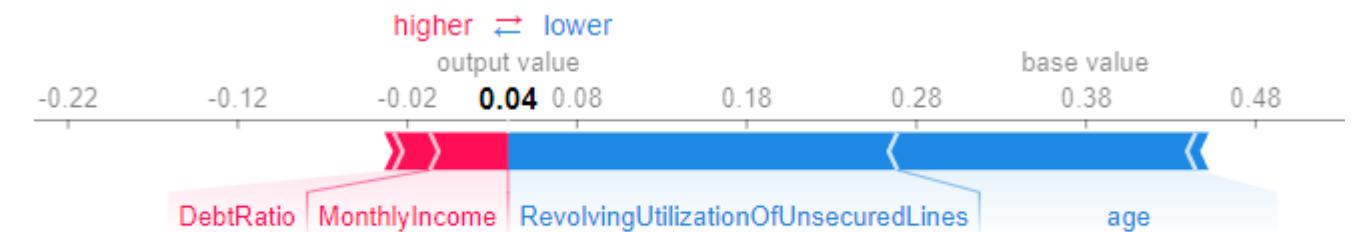


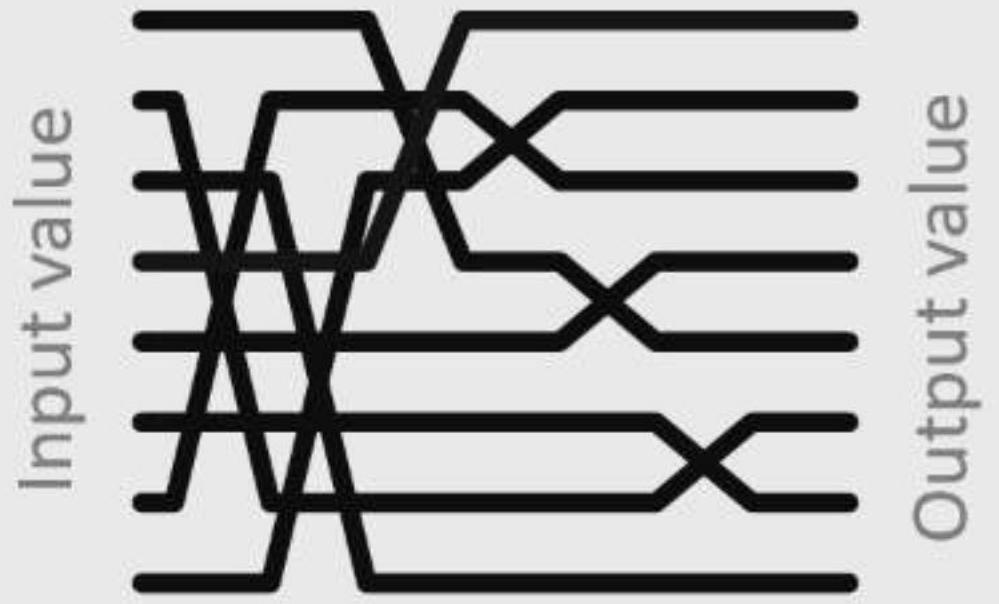
# Types Of Explanations

Global  
(Model-wide)

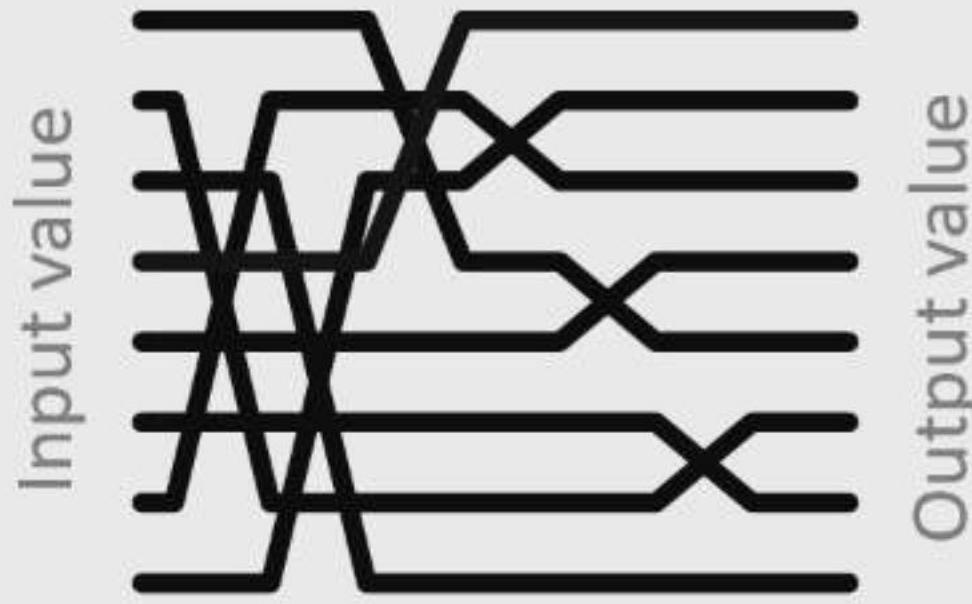


Local  
(Single Prediction)

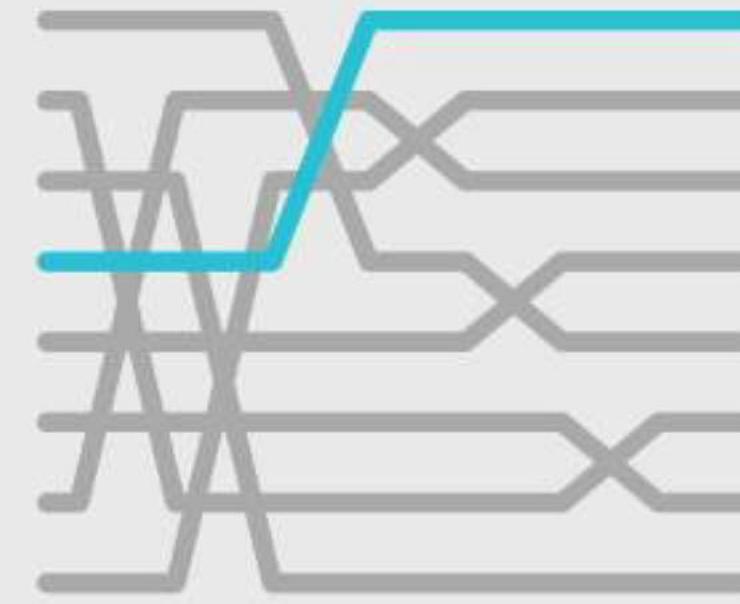




Complex models are  
inherently complex!



Complex models are  
inherently complex!



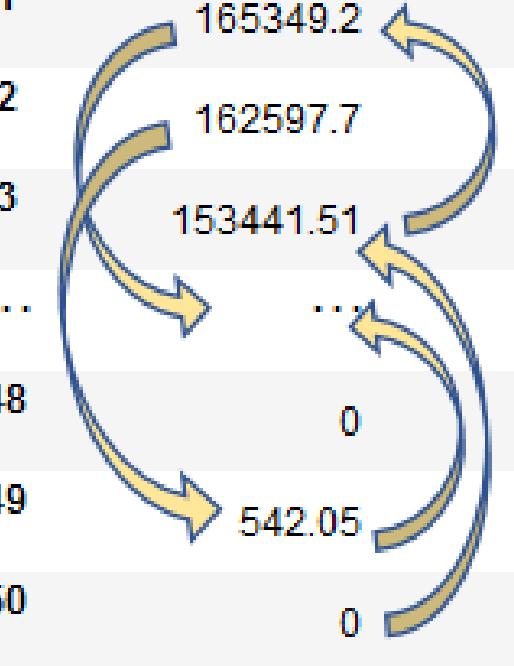
But a single prediction involves only a  
small piece of that complexity.

# Permutation Importance

# Permutation Importance

- Global, Black Box
- Measure model performance decreases when feature is not available (permuted)

	RD Spend	Administration	Marketing Spend	Profit	state_California
1	165349.2	136897.8	471784.1	192261.83	0
2	162597.7	151377.59	443898.53	191792.06	1
3	153441.51	101145.55	407934.54	191050.39	1
...	...	...	...	...	...
48	0	135426.92	0	42559.73	1
49	542.05	51743.15	0	35673.41	0
50	0	116983.8	45173.06	14681.4	1





Search projects



Help

Donate

Log in

Register

# eli5 0.10.1



Latest version

pip install eli5

Last released: Aug 29, 2019

Debug machine learning classifiers and explain their predictions

## Navigation

Project description

Release history

Download files

## Project links

Homepage

## Statistics

## Project description

pypi v0.10.1 build passing codecov 97% docs failing

ELI5 is a Python package which helps to debug machine learning classifiers and explain their predictions.

hi there, i am here looking for some help. my friend is a interic graphics software on pc. any suggestion on which software to sophisticated software(the more features it has,the better)



It provides support for the following machine learning frameworks and packages:

# Demo

Permutation Importance with eli5

# Permutation Importance - Advantages

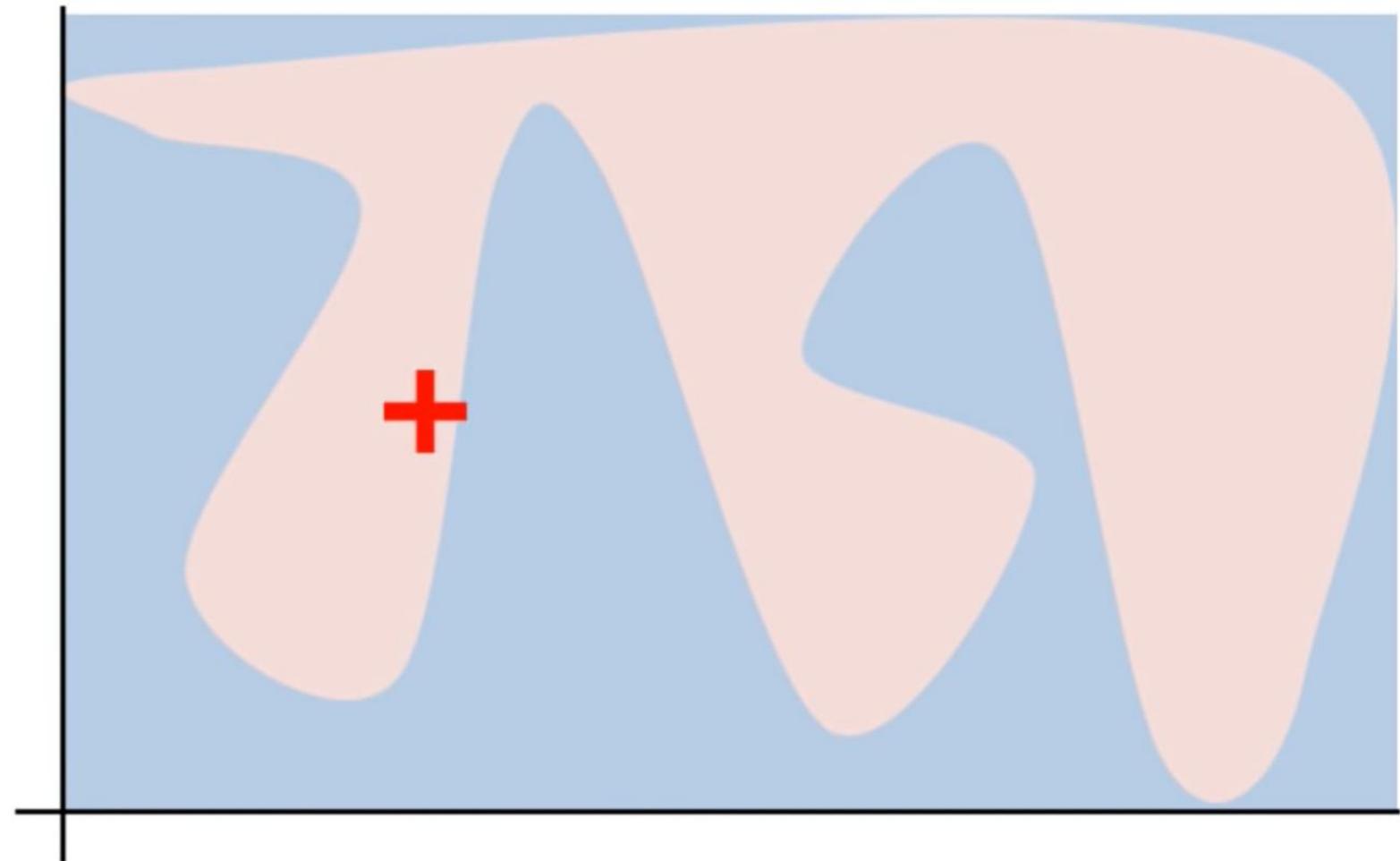
- Global insight into the model
- Useful for feature selection
- Very fast

# Permutation Importance - Disadvantages

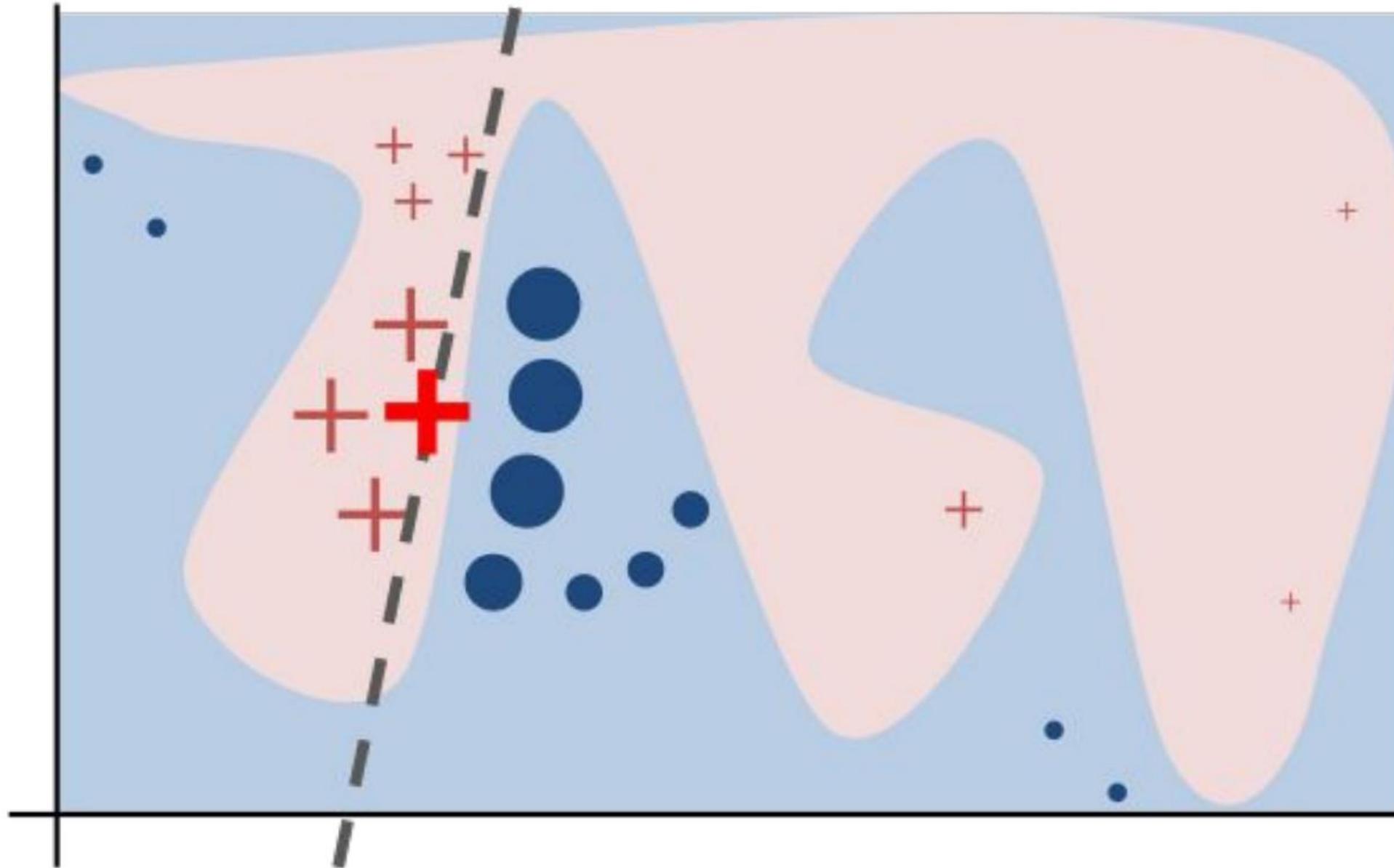
- Correlated features are a problem
- Not clear whether to use training or testing data

LIME

**L**OCAL  
**I**NTERPRETABLE  
**M**ODEL-AGNOSTIC  
**E**XPLANATIONS

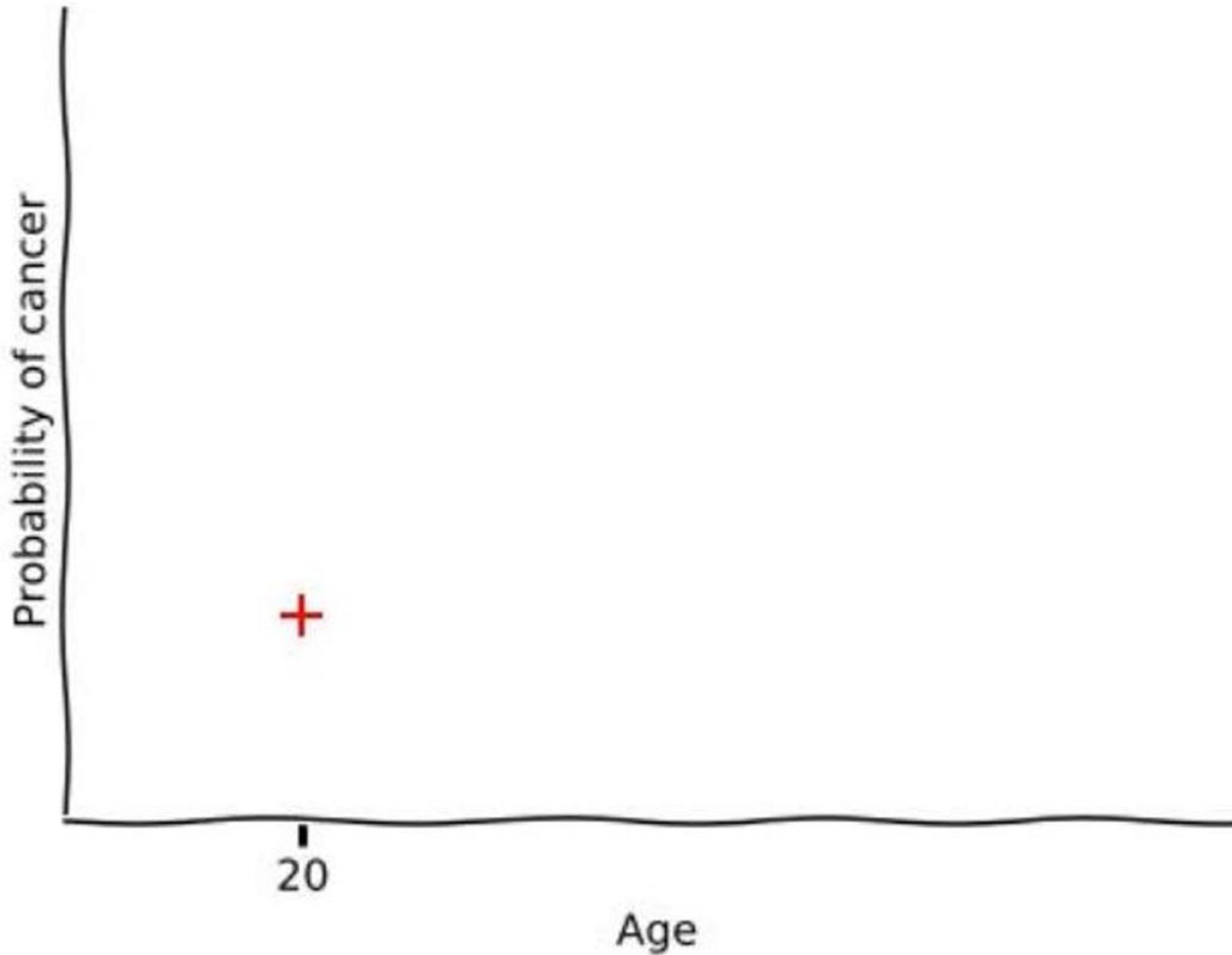


# How LIME Works

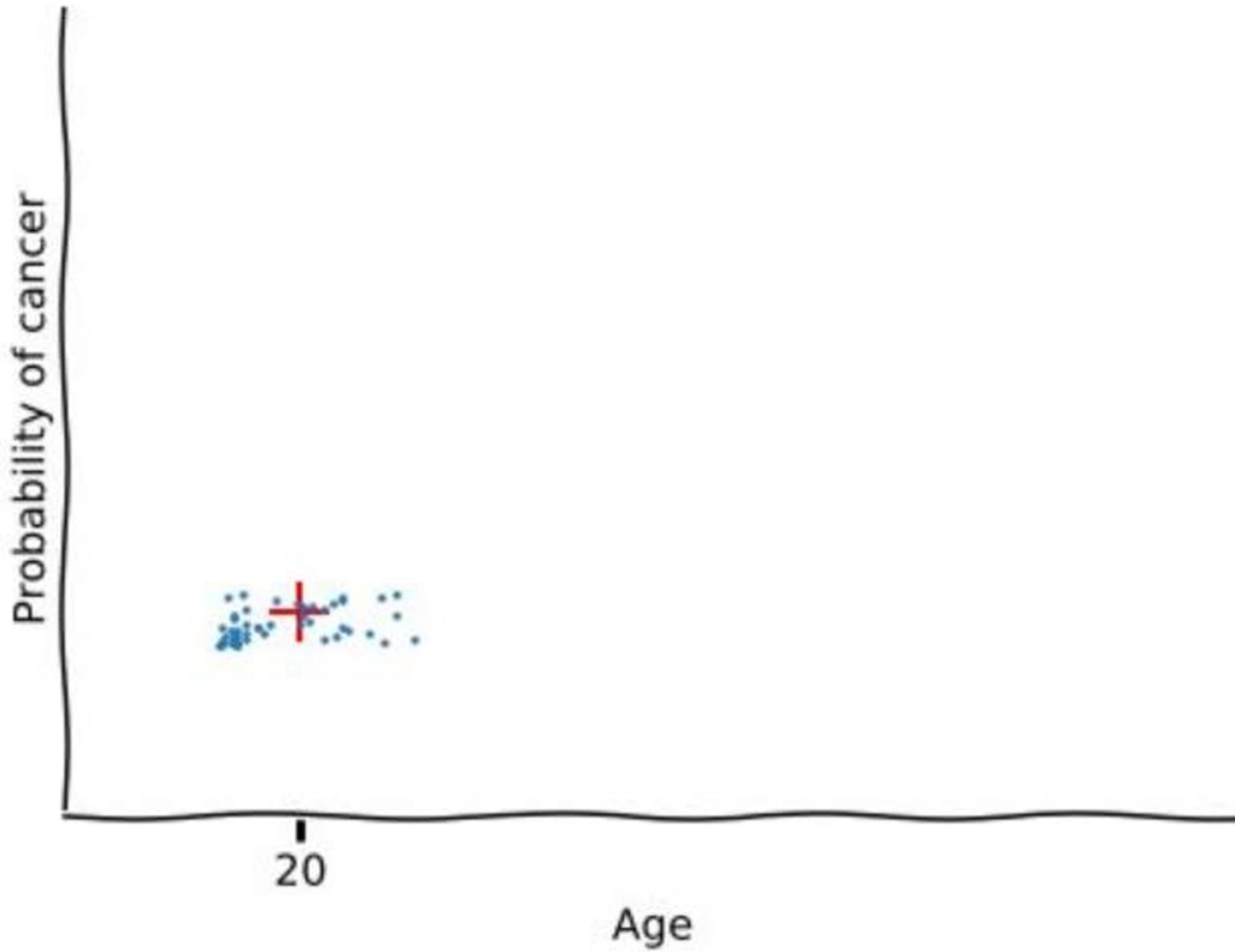


"Why Should I Trust You?": Explaining the Predictions of Any Classifier  
*Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin*

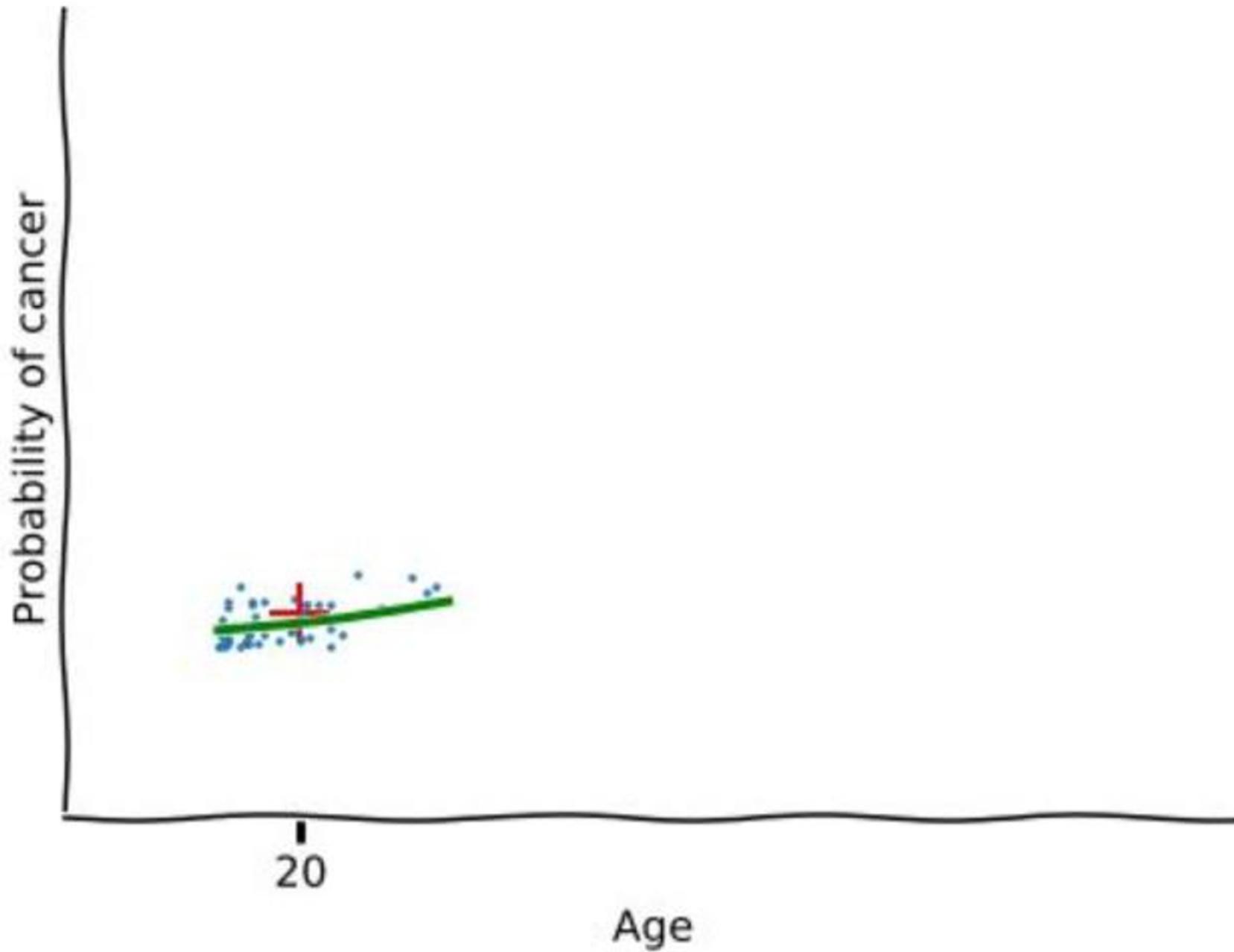
# How LIME Works



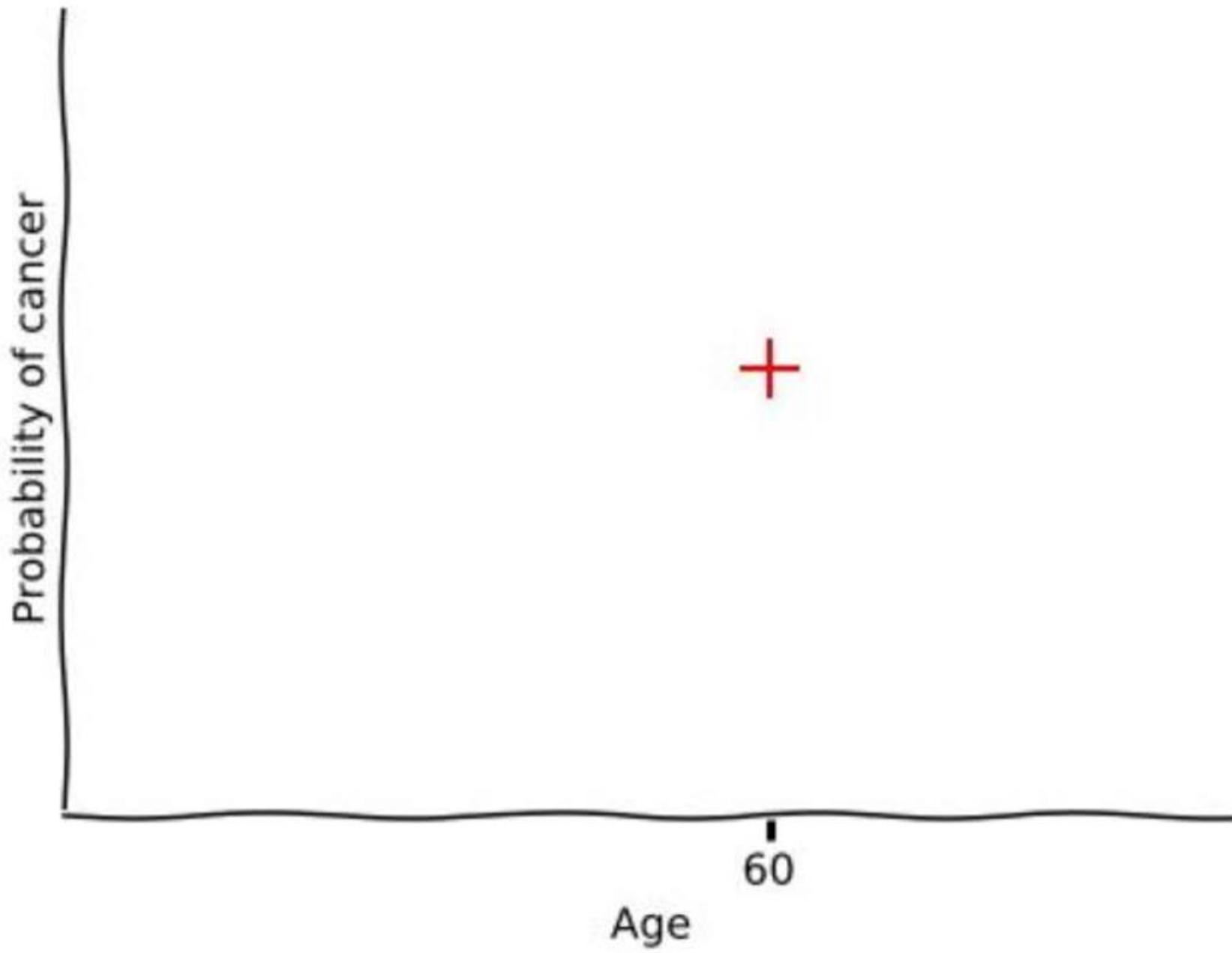
# How LIME Works



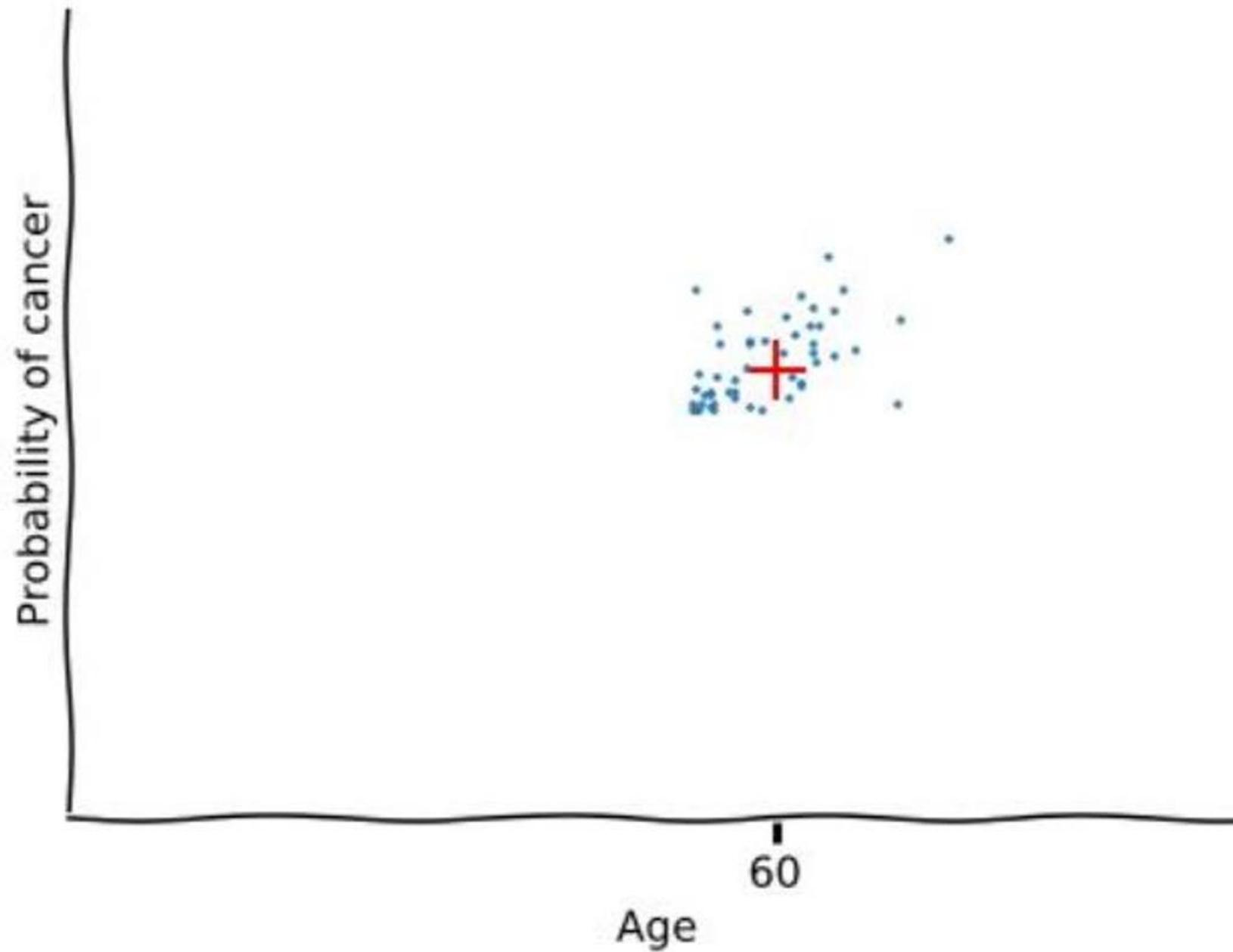
# How LIME Works



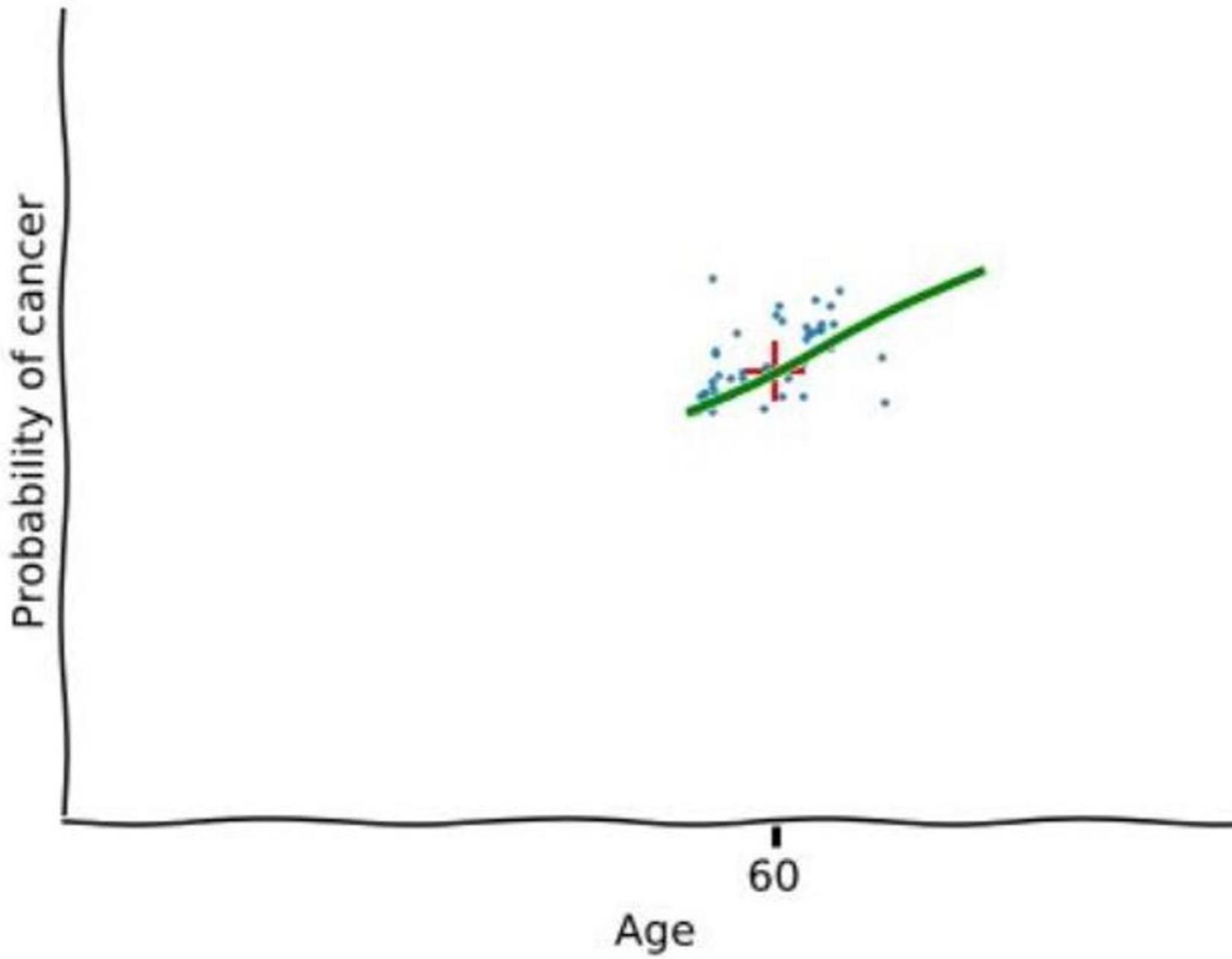
# How LIME Works



# How LIME Works



# How LIME Works



# How LIME Works

1. Choose an **observation** to explain
2. Create new dataset around observation by sampling from training data distribution
3. Calculate distances between new points and observation, that's our measure of **similarity**
4. Use model to predict class of the new points
5. Find the subset of **m** features that has the strongest relationship with our target class
6. Fit a linear model on fake data in **m** dimensions weighted by similarity
7. Weights of linear model are used as explanation of decision

## LIME Also Works With Image Data



Original Image



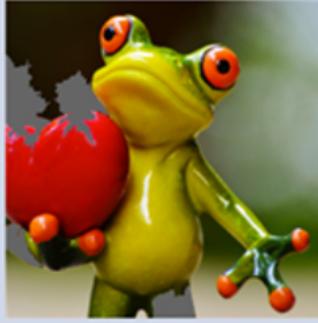
Interpretable  
Components

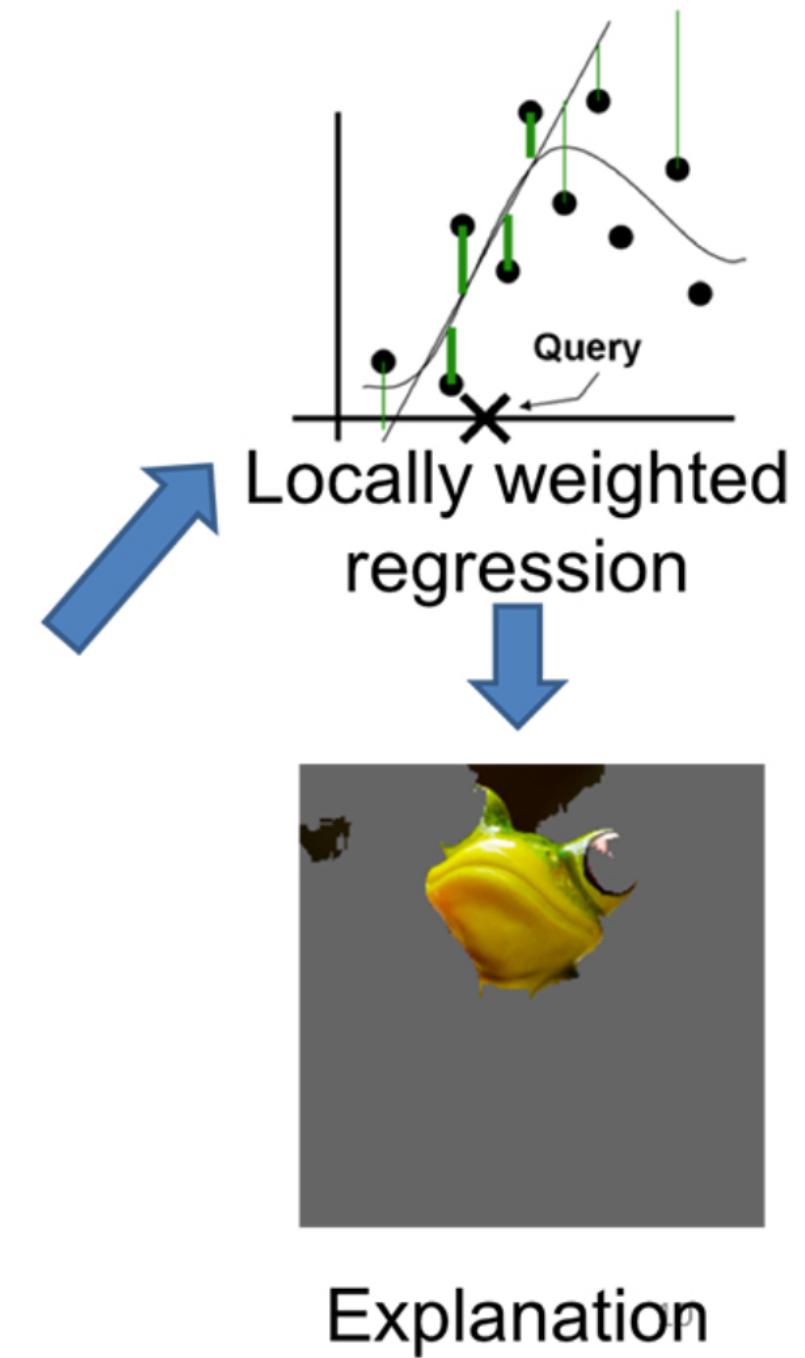
# LIME Also Works With Image Data



Original Image  
 $P(\text{tree frog}) = 0.54$



Perturbed Instances	$P(\text{tree frog})$
	0.85
	0.00001
	0.52

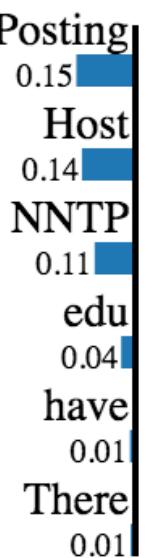


# LIME Also Works With Text Data

Prediction probabilities



atheism



christian

## Text with highlighted words

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.

# Demo

Local Explanations with LIME

# LIME - Advantages

- Local explanations, useful for end users
- Works with tabular data, images and text
- Black box

# LIME - Disadvantages

- Depends on random sampling, can be unstable
- Linear model surrogate can be inaccurate
  - This can be tested with R-squared score
- Can be slow for complex models

# SHAP

# SHAP

Lundberg and Lee. A unified approach to interpreting model predictions  
NeurIPS 2017 (*oral presentation*)

Lundberg and Lee. An unexpected unity among methods for interpreting model predictions  
NeurIPS Workshop on Interpretable Machine Learning in Complex Systems 2016 (*best paper award*)



0  
↓



Base rate

20%

0

$E[f(x)]$





Base rate

20%

0



$E[f(x)]$



Prediction for John

55%

$f(x)$





Base rate

20%

0

$E[f(x)]$

Prediction for John

55%

$f(x)$

→

How did we get here?



20%

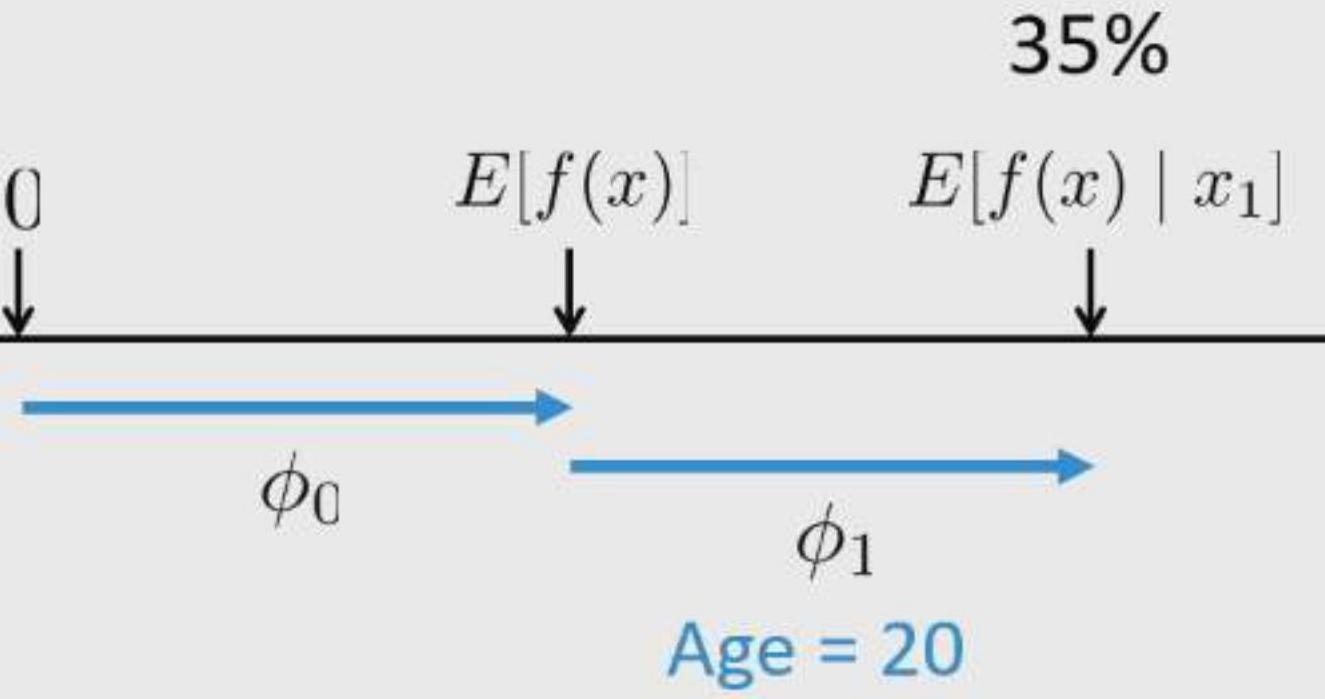
0

$E[f(x)]$



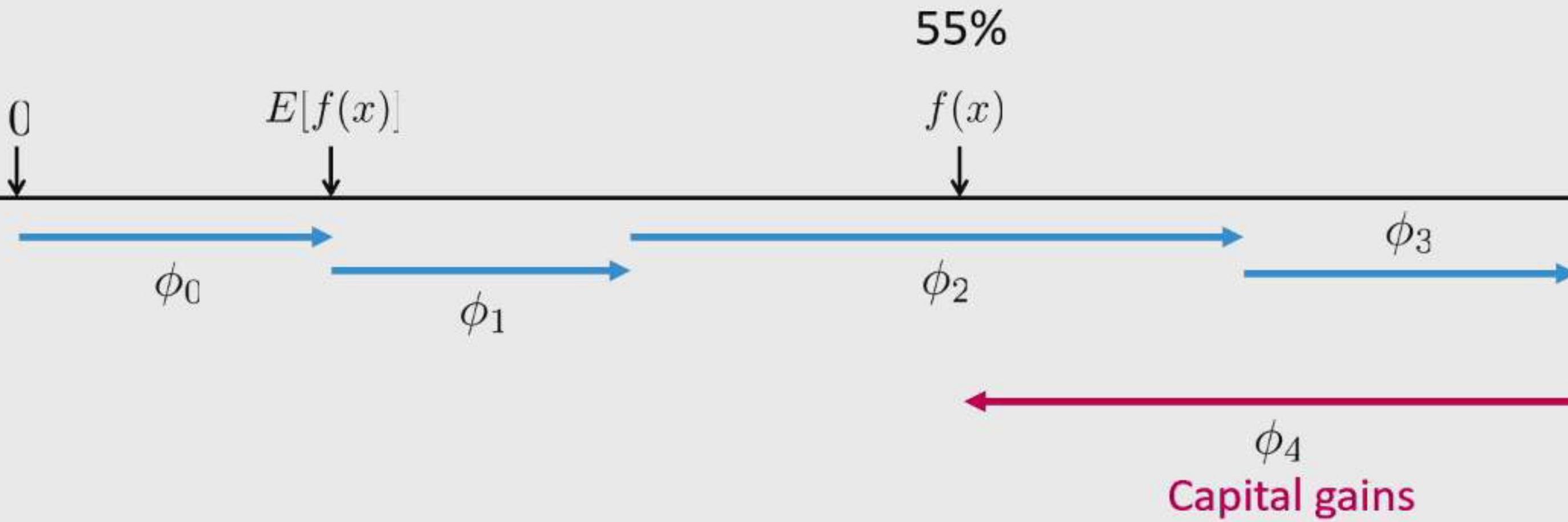
$\phi_0$

Base rate

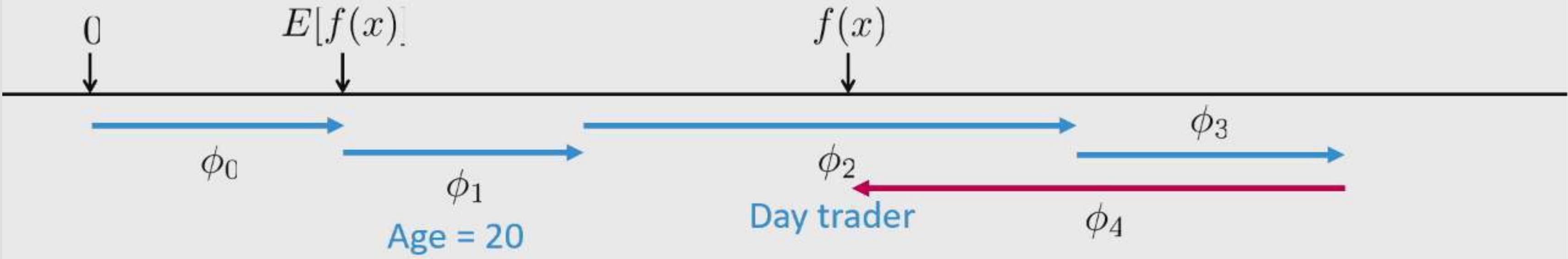




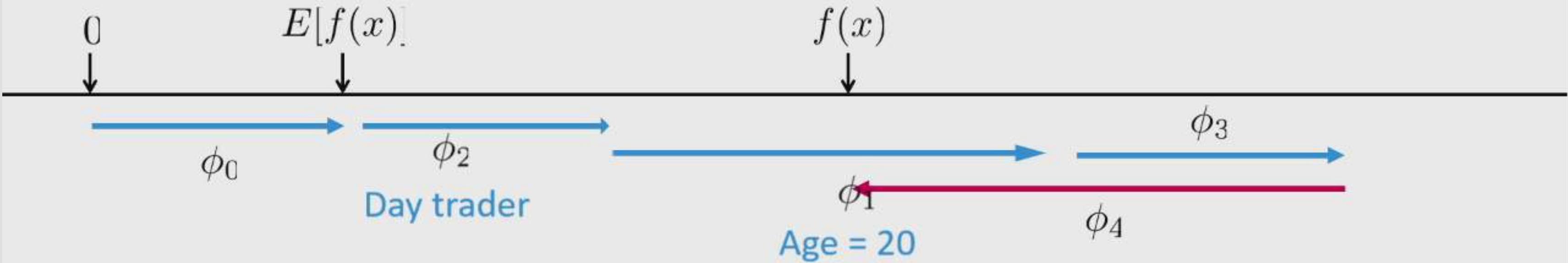




**The order matters!**



**The order matters!**

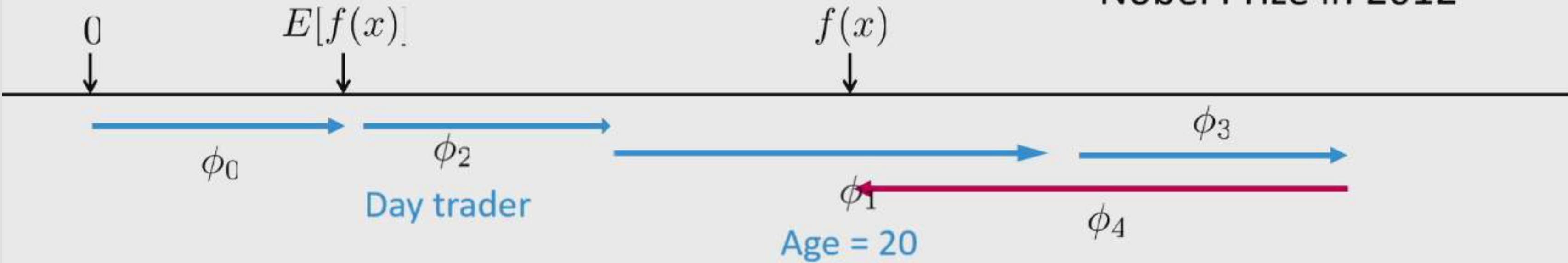


Lloyd Shapley

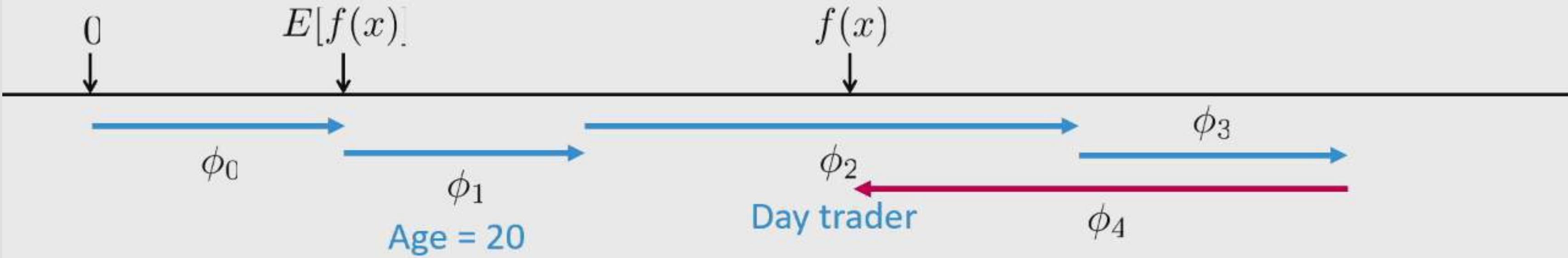


The order matters!

Nobel Prize in 2012

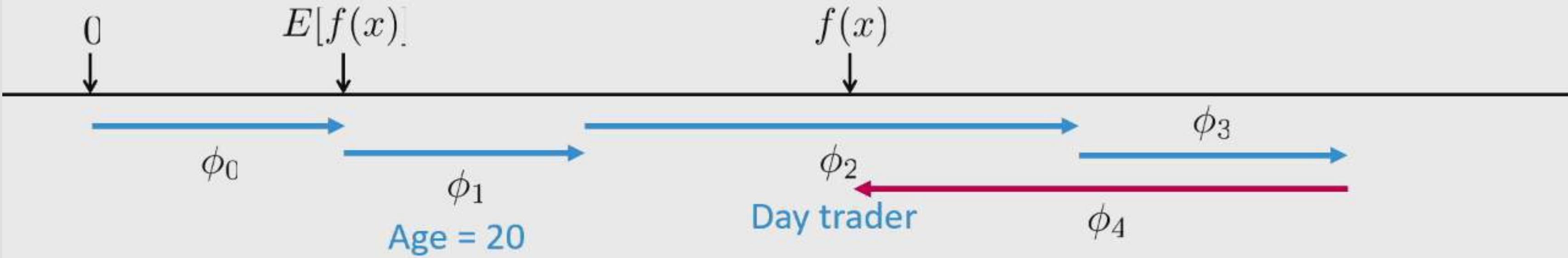


Shapley values result from **averaging over all  $N!$  possible orderings**.



# SHapley Additive exPlanation (SHAP) values

Shapley values result from **averaging over all  $N!$  possible orderings**.



# How SHAP Works

The weight of each feature is computed using the Shapley values method from game theory.

To get the importance of feature  $X\{i\}$ :

- Get all subsets of features  $S$  that do not contain  $X\{i\}$
- Compute the effect on our predictions of adding  $X\{i\}$  to all those subsets

That can be computationally expensive, but SHAP has optimizations for different models (linear, trees, etc..)

# Demo

Local And Global Explanations with shap

# SHAP - Advantages

- Solid theoretical foundation in game theory
- Generates both local and global explanations
- Black box
  - optimizations are available for Glass box models

# SHAP - Disadvantages

- Ignores feature dependence
- Black box (kernel) SHAP is slow

LIME

DeepLIFT

Shapley reg. values

SHAP

Relevance prop.

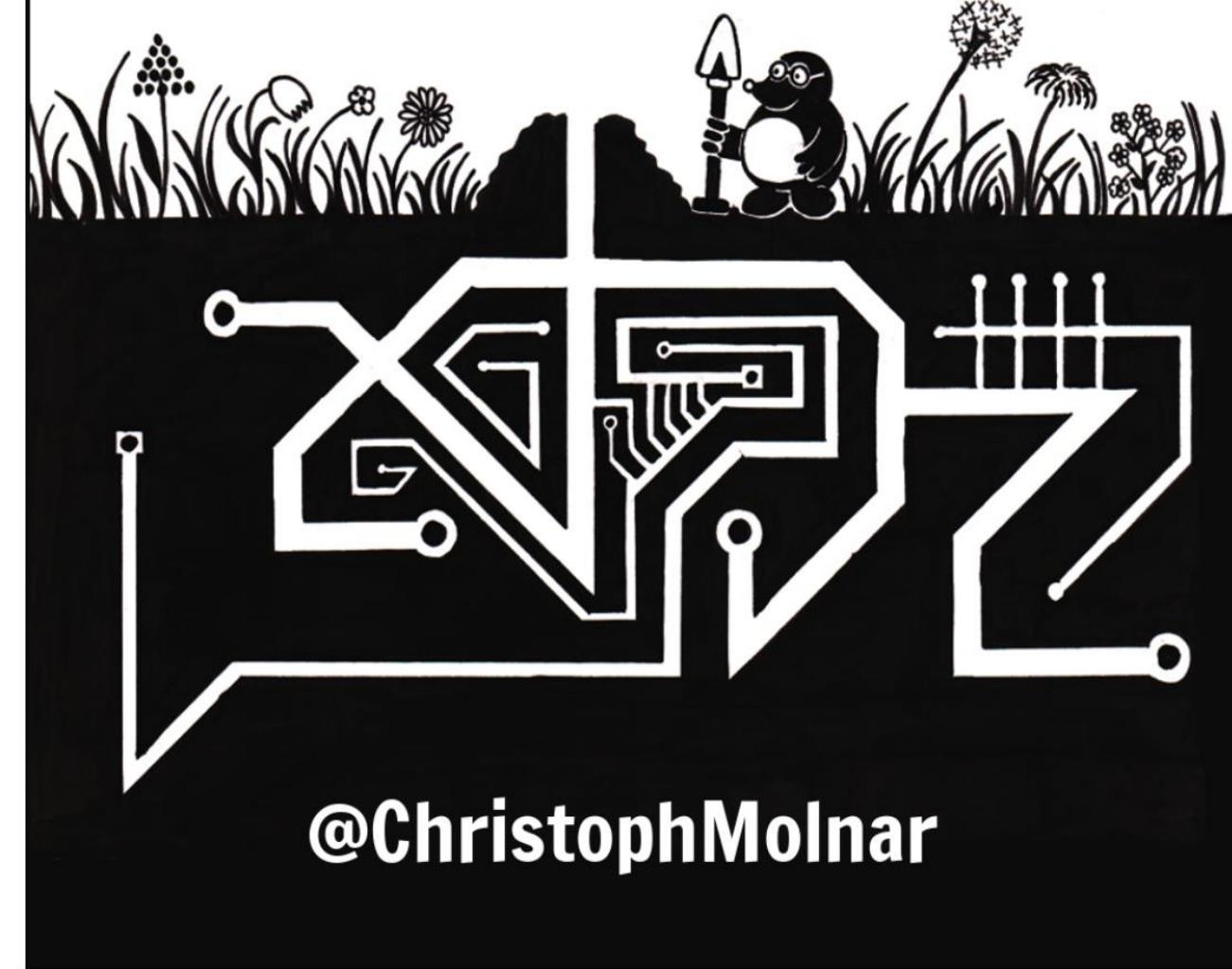
QII

Shapley sampling

Saabas

# Interpretable Machine Learning

A Guide for Making  
Black Box Models Explainable



@ChristophMolnar

# Questions?

**CODECAMP** ❤ **FEEDBACK**



[codecamp.ro/feedback](http://codecamp.ro/feedback)

# Thank You!