

MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME

ABOUT ME

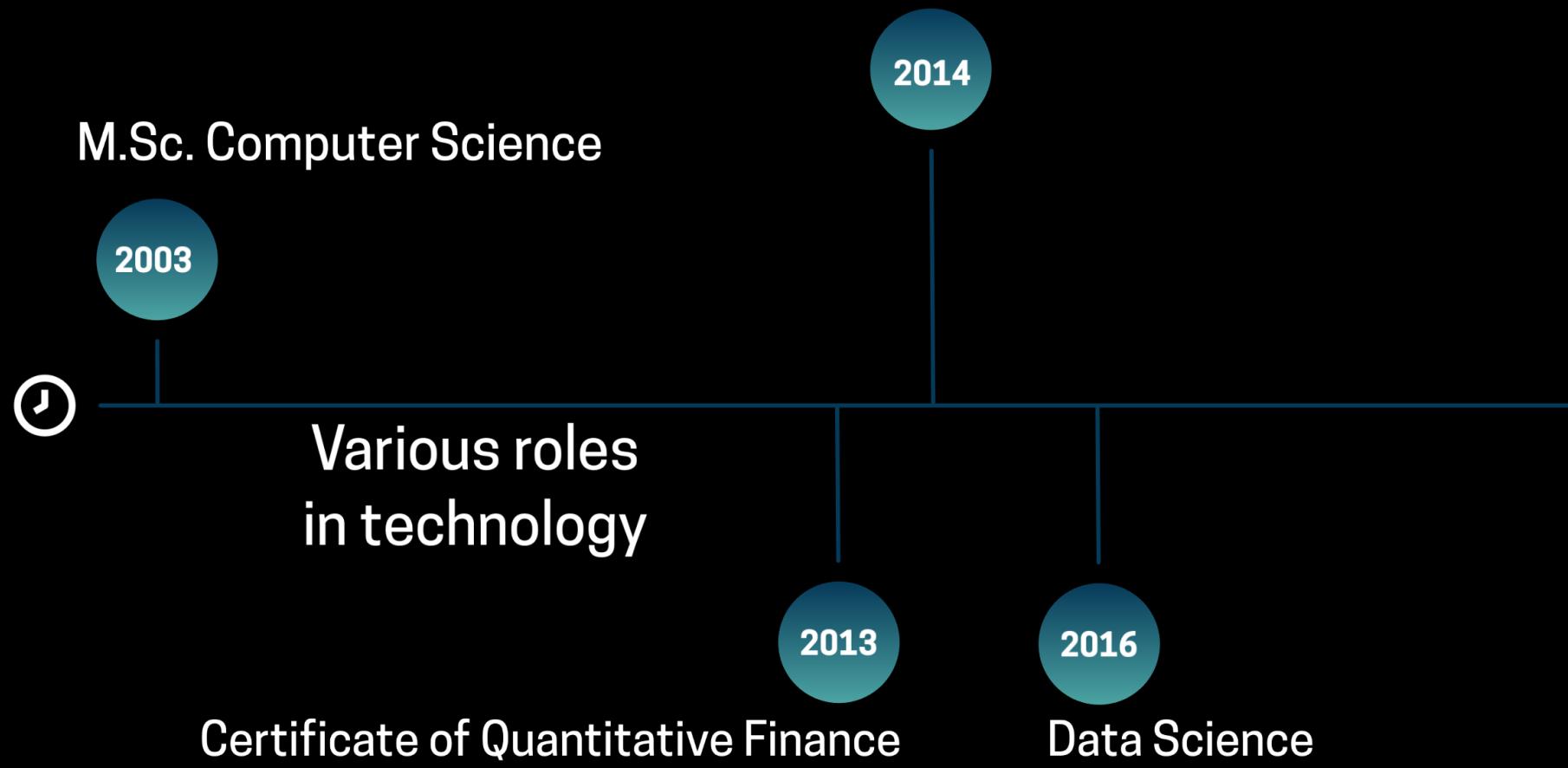
RESUME

KEY
INTERESTS

HOBBIES

RESUME

Quantitative analyst
(Market Risk)

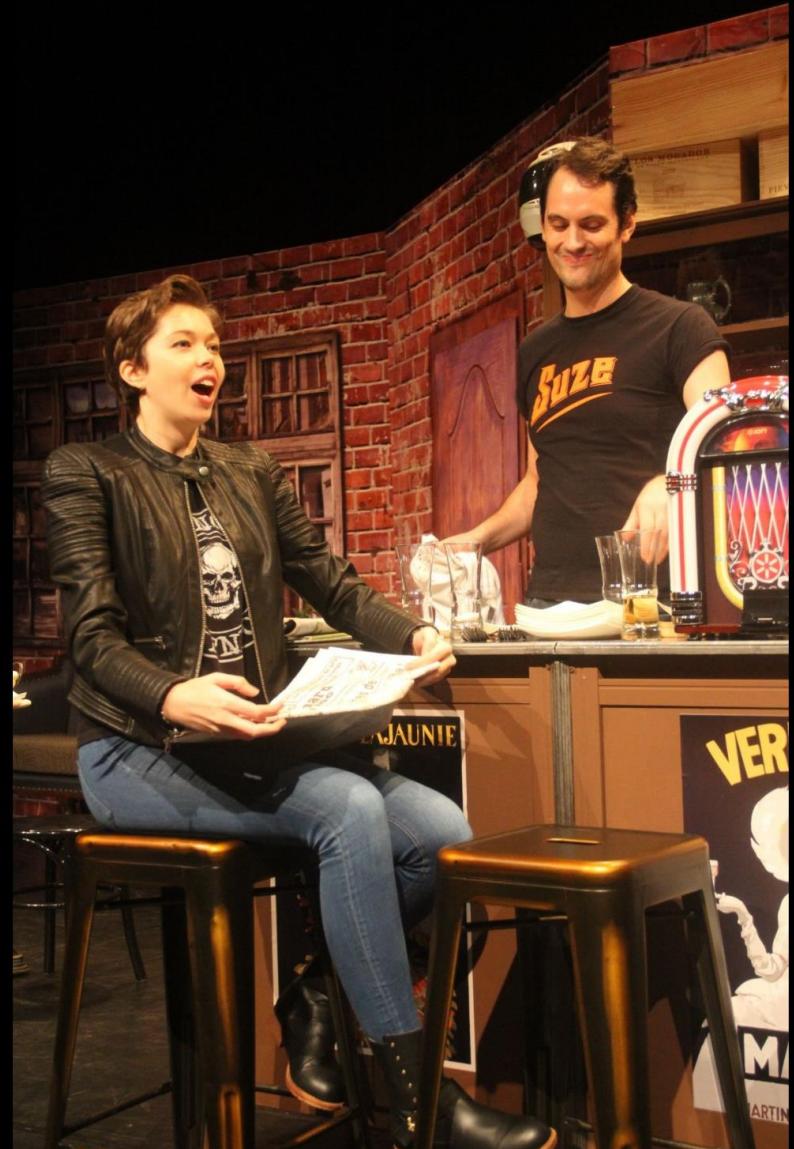


KEY INTERESTS

- Quantitative investment
- Systematic trading
- Machine learning
- Blockchain
- System development

HOBBIES

- Traveling
- Acting (Theater)
- Squash
- Cooking



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME



INTRO

MODEL

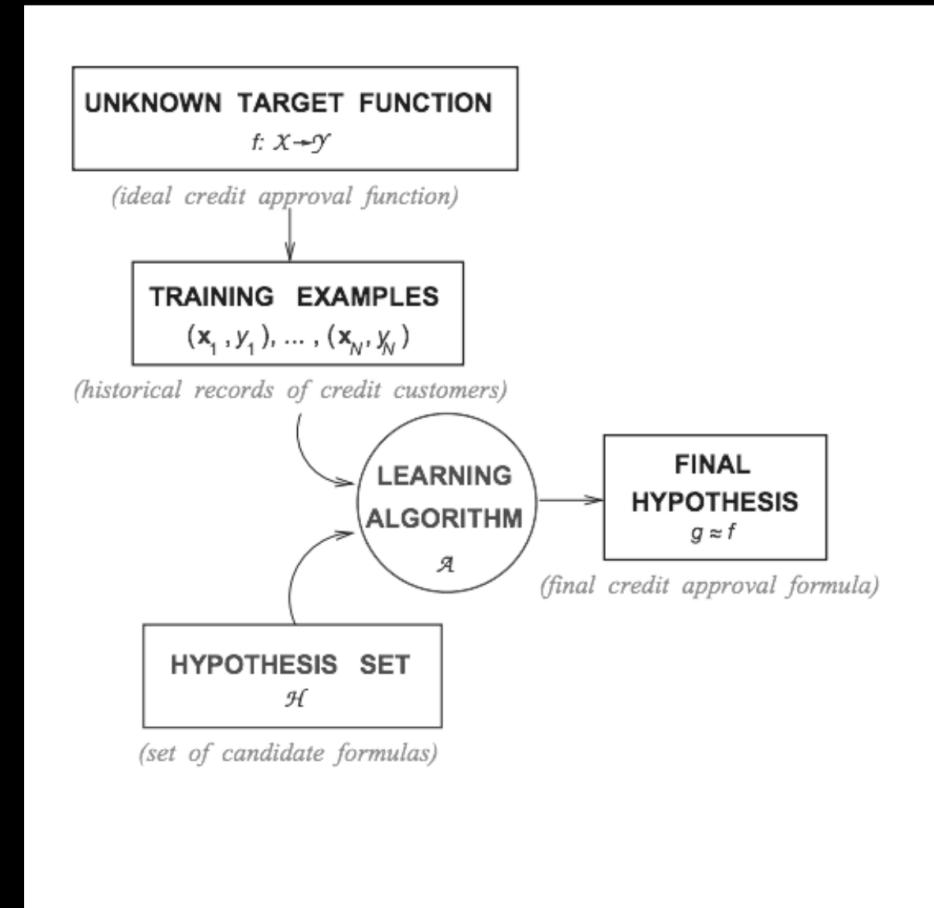
TECHNIQUES

AI IN
FINANCE

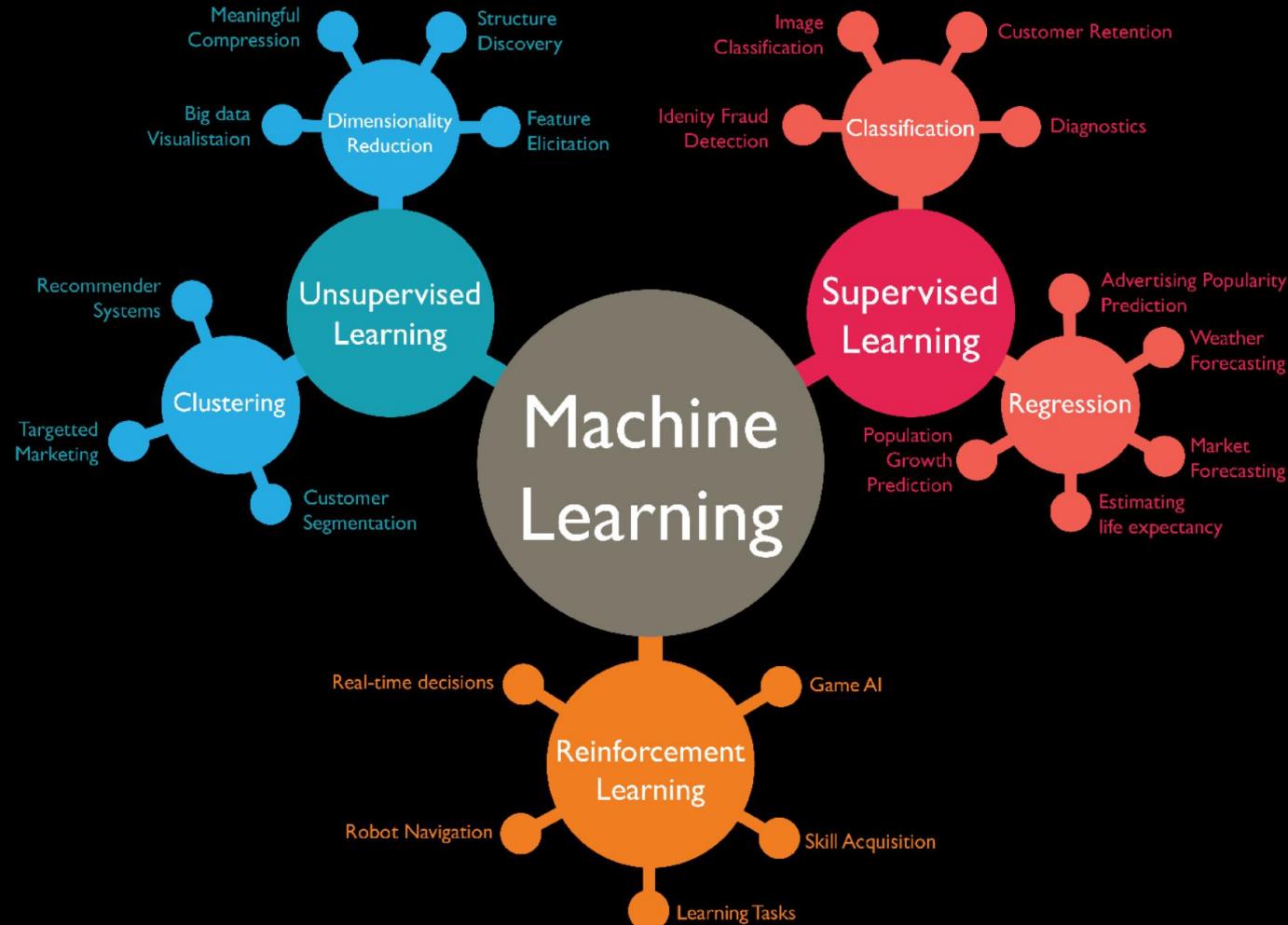


WHAT IS A MODEL?

Machine learning (ML) models are algorithm that allows software applications to become more accurate in predicting outcomes without being explicitly programmed. The basic premise of machine learning is to build algorithms that can be trained with data and use statistical analysis to predict an outcome.

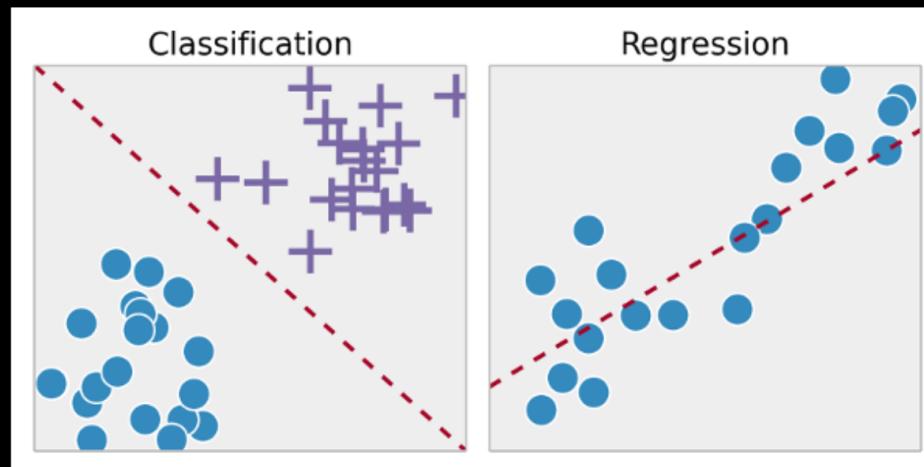


TECHNIQUES

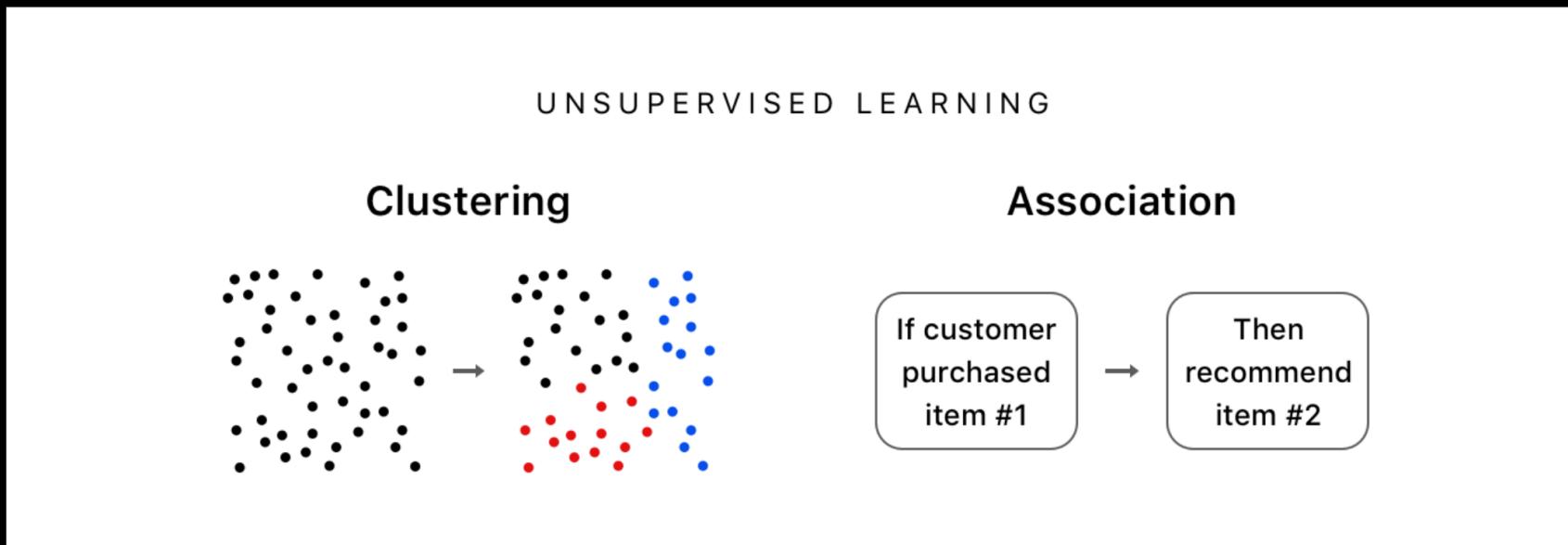


SUPERVISED LEARNING OUTCOME

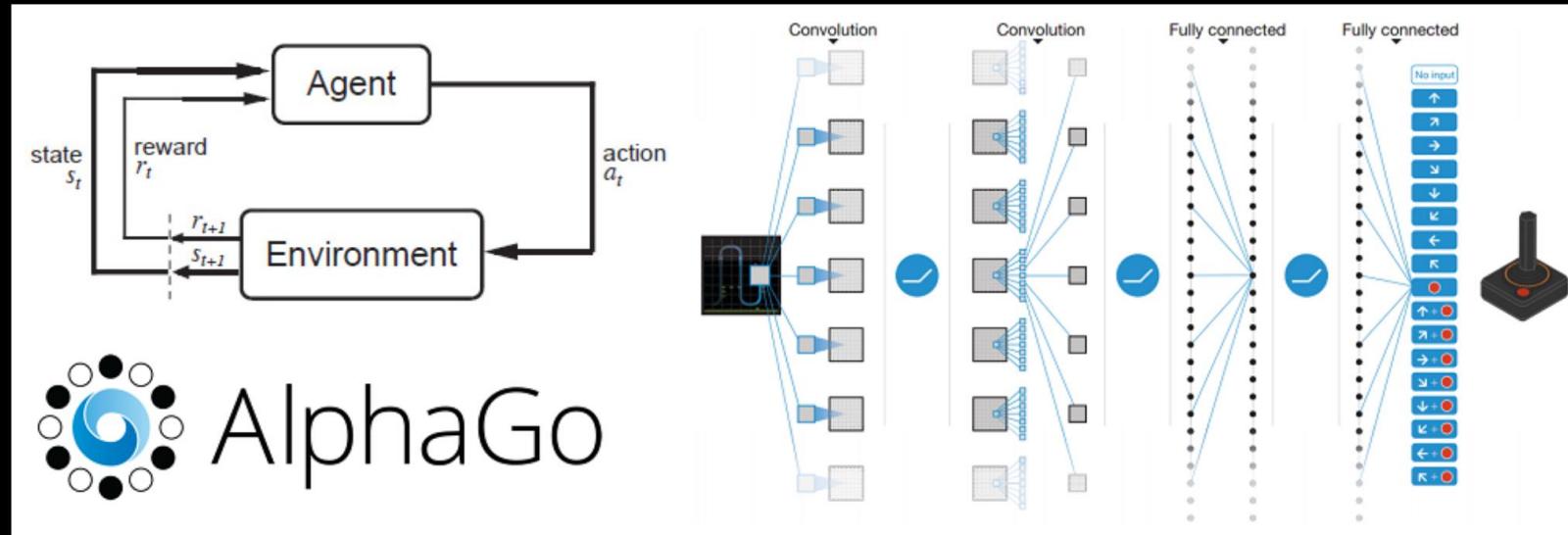
- A classification is the problem of predicting a discrete class label output. For an example, predict if the stock price of Microsoft will go up or down tomorrow. Result is often expressed as probabilities (e.g. 10% down, 90% up). Accuracy = correct predictions / total predictions.
- A regression is the problem of predicting a continuous quantity output. For example, predict the stock price of Microsoft tomorrow. Many ways to estimate accuracy. RMSE is probably most common



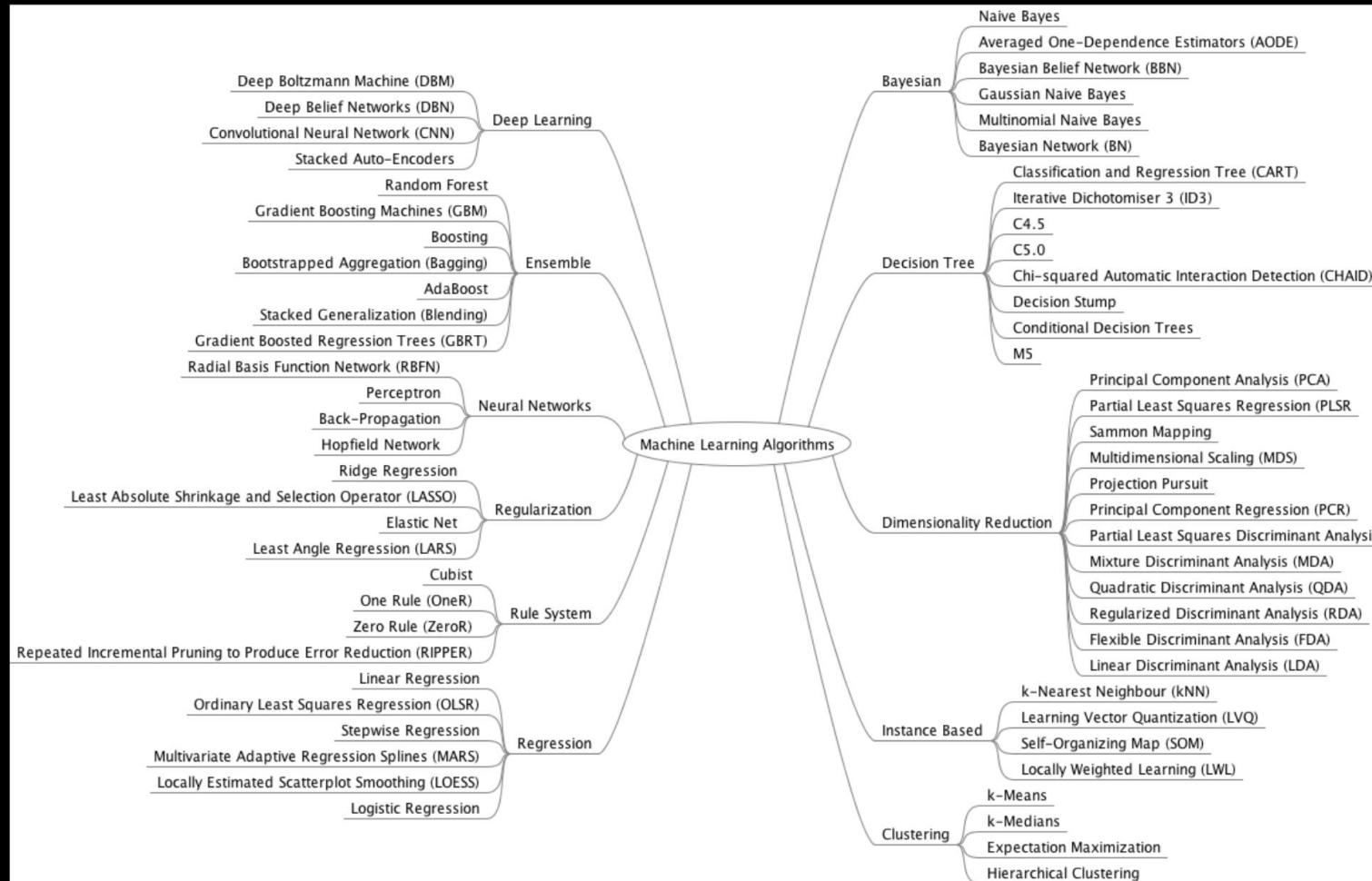
UNSUPERVISED LEARNING



REINFORCEMENT LEARNING



TECHNIQUES



ARTIFICIAL INTELLIGENCE IN FINANCE

Marketing

- Client profiling
- Recommendation of financial products

Process automation

- Chatbots
- Call-center automation
- Paperwork automation (e.g. Damage estimation for insurance)
- Gamification of employee training, and more.
- Virtual recruiter

Security

- Financial monitoring (Money laundering, etc.)
- Network security

Underwriting and credit scoring

- Credit scoring
- Current account / Mortgage portfolio modeling

Algorithmic trading

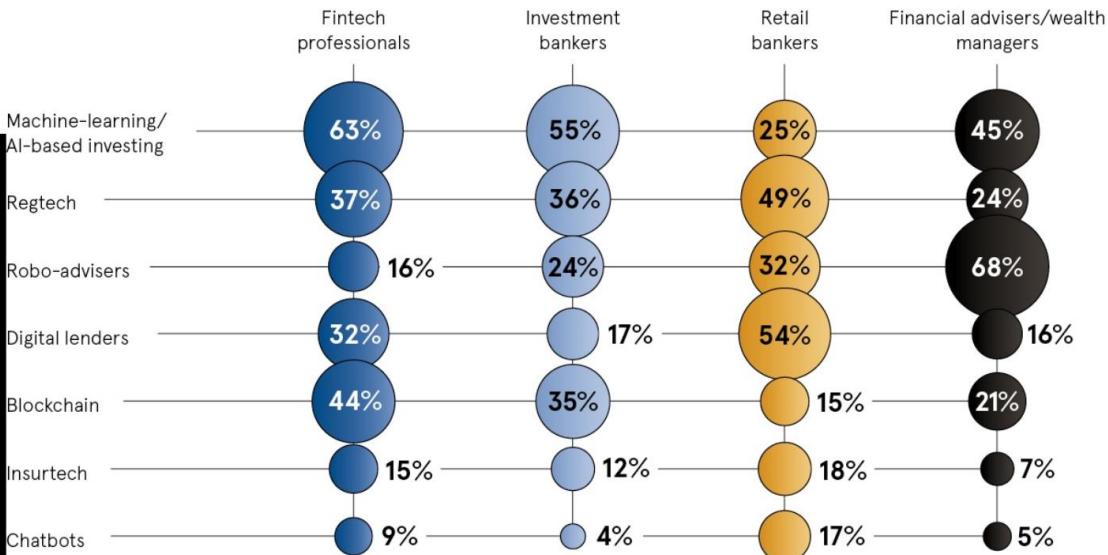
- AI trading
- Market prediction with alternative data (satelite)

Asset Management

- Robo-Advisor / ETFs /
- Sentiment Analysis

Most important technologies disrupting the financial world

Percentage of different industries who believe the following are important to their sector



LinkedIn 2017



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

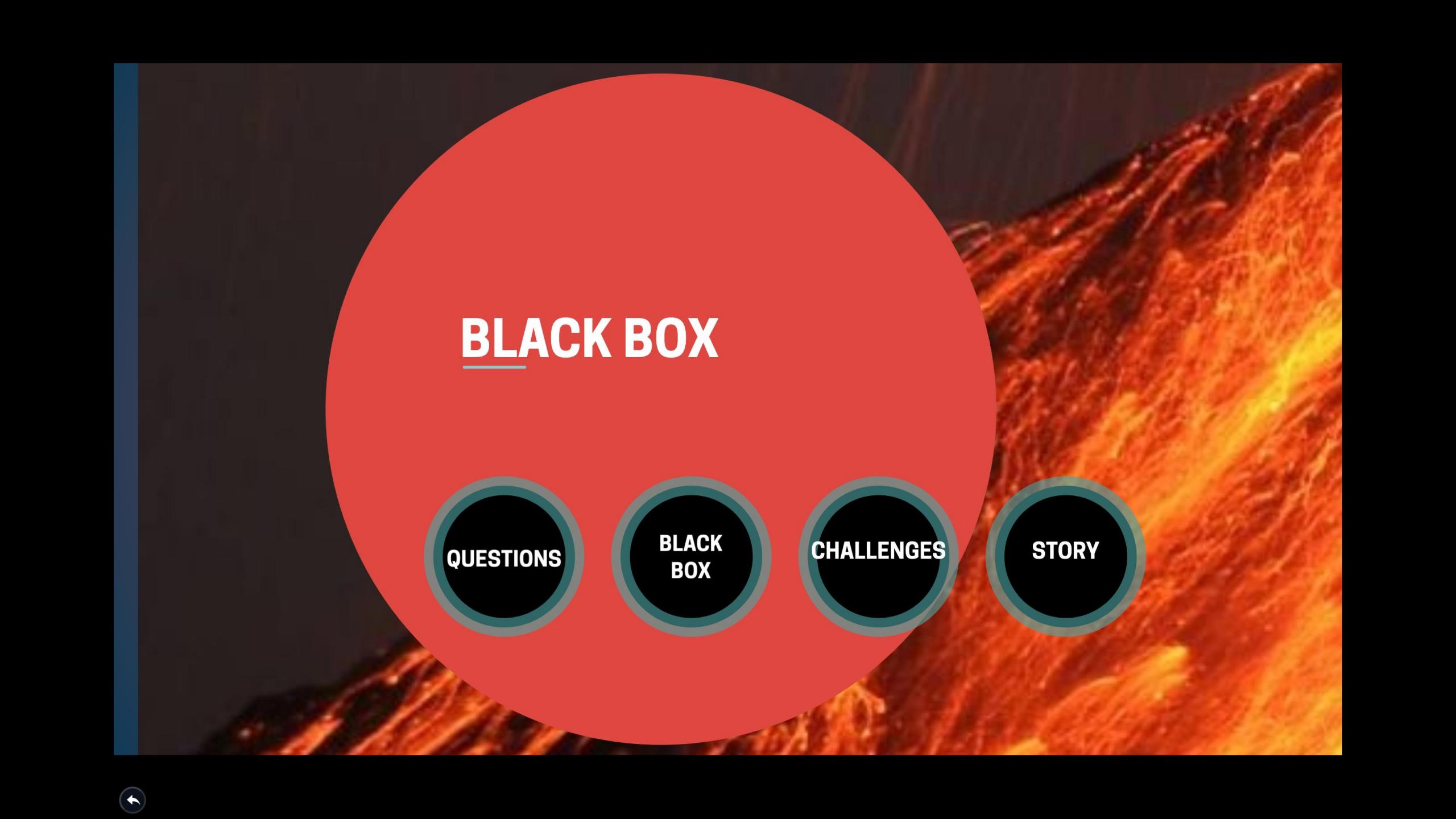
BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME



BLACK BOX

QUESTIONS

BLACK
BOX

CHALLENGES

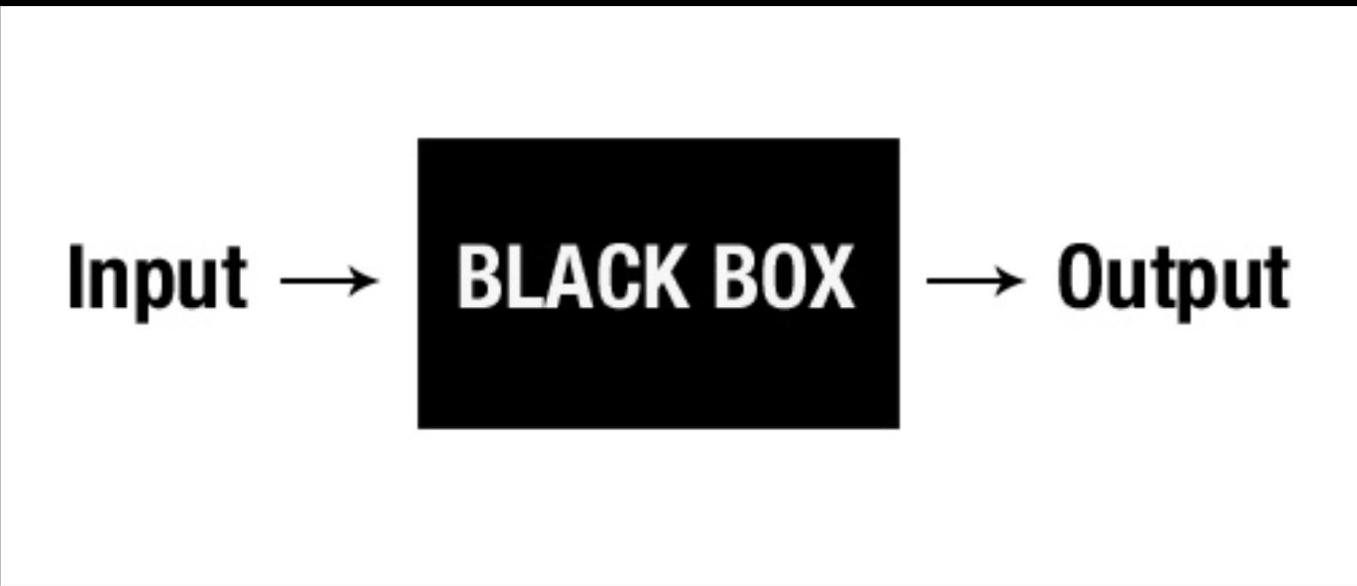
STORY

QUESTIONS WE NEED TO ANSWER

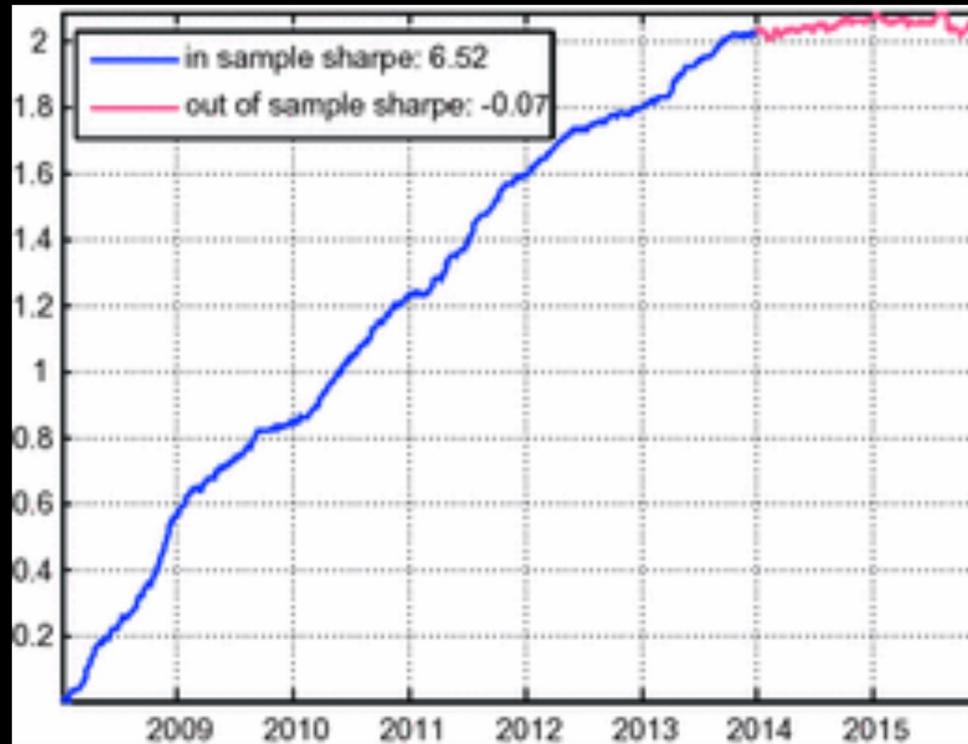
- Should I offer a loan to this customer?
- What should be the insurance premium for this customer?
- Is this transaction related to terrorism financing?



BLACK BOX



BLACK BOX



CHALLENGES

- Is this model robust?
- Do i understand this model?
- Can i trust this model?
- How can i explain this model?

STORY

Education

‘Creative ... motivating’ and fired



Sarah Wysocki was out of work for only a few days after she was fired by DCPS last year. She is now teaching at Hybla Valley Elementary School in Fairfax County. (Jahi Chikwendiu/The Washington Post)

By [Bill Turque](#)

March 6, 2012



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME

HOW DID WE GET THERE?



MODEL

SIMPLE
MODEL

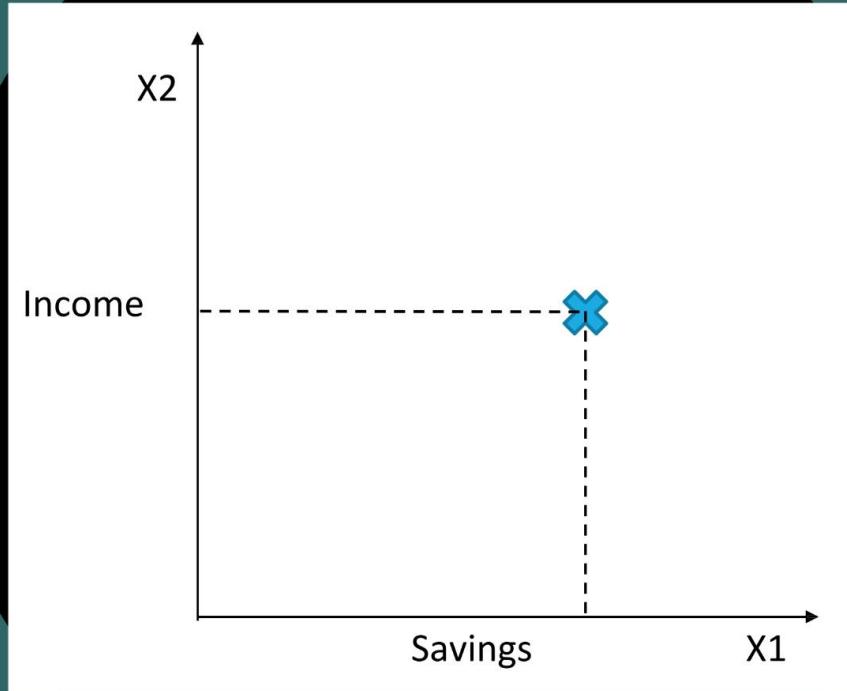
TREE

MORE
COMPLEX

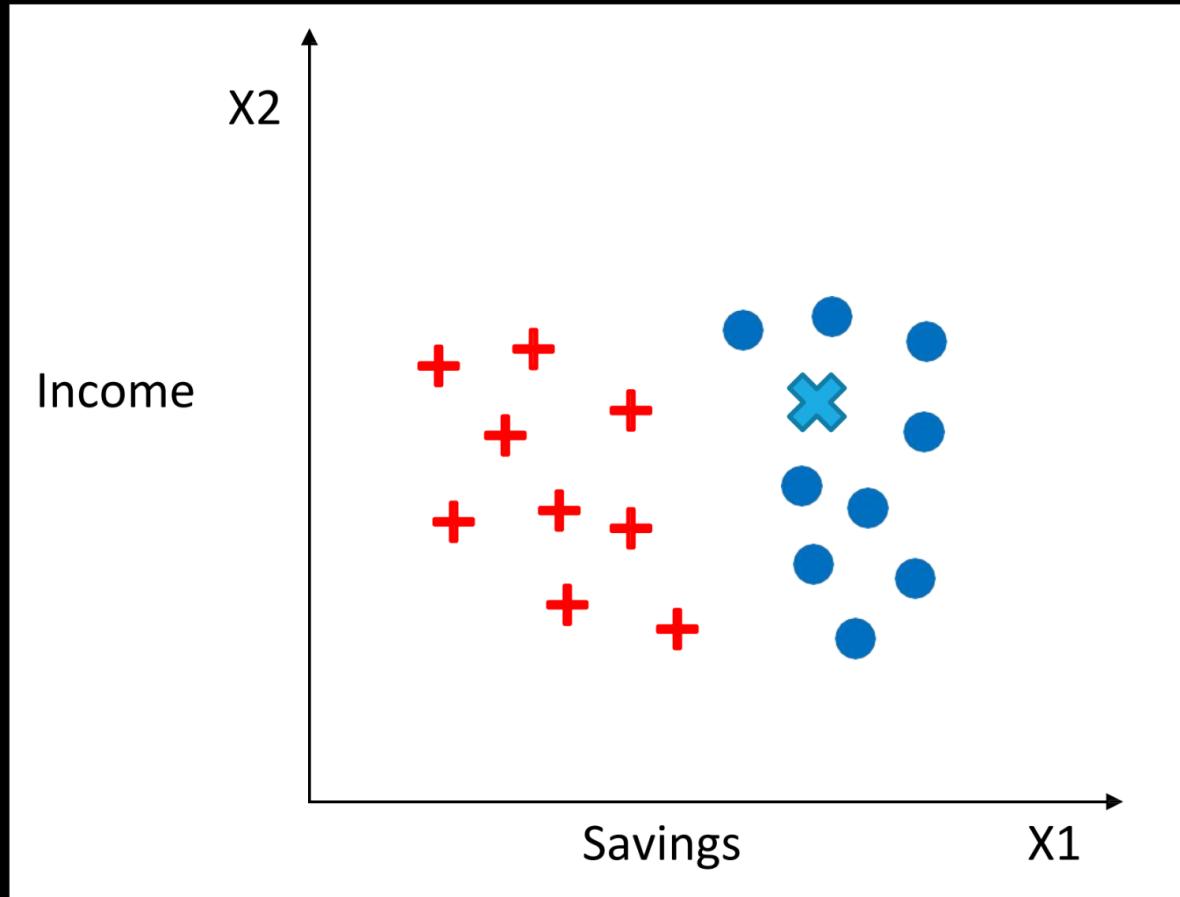
ACCURACY /
INTERPRETABILITY

NEXT?

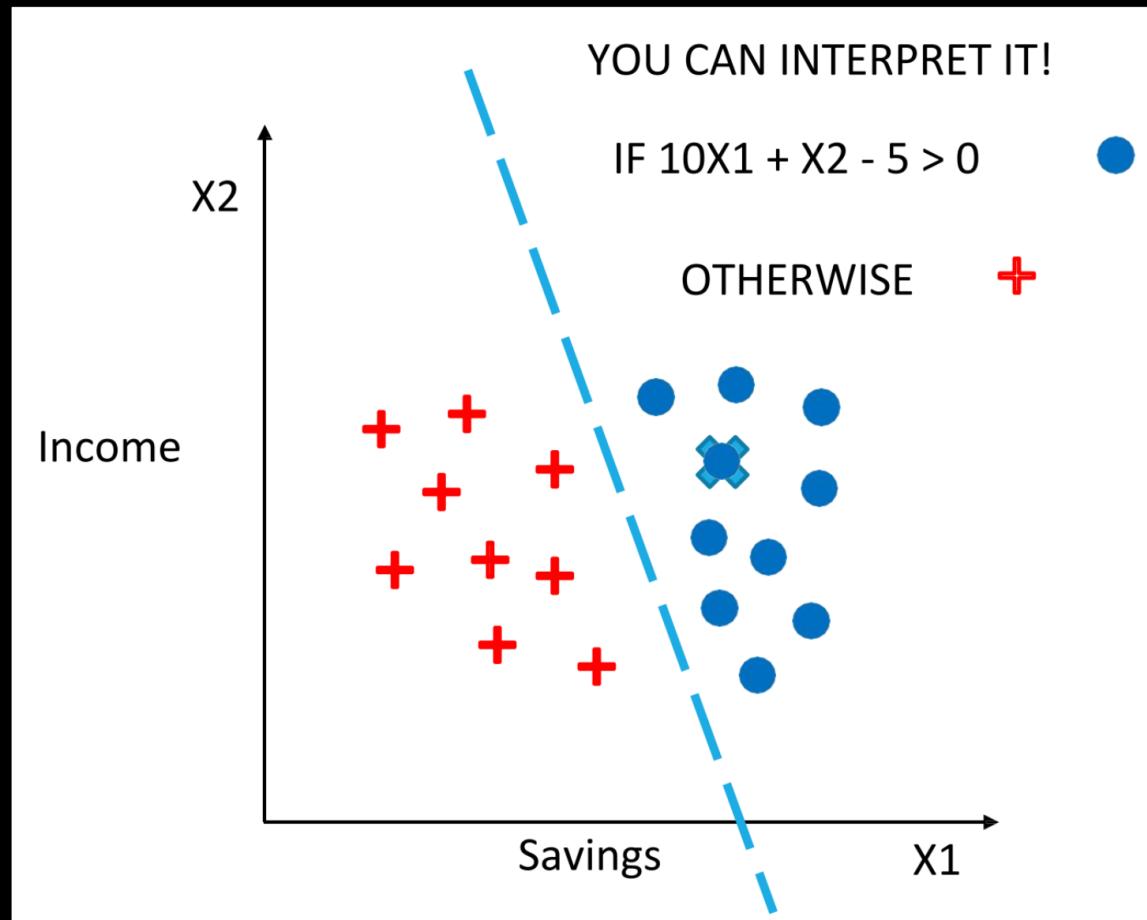
WILL THE LOAN DEFAULT?



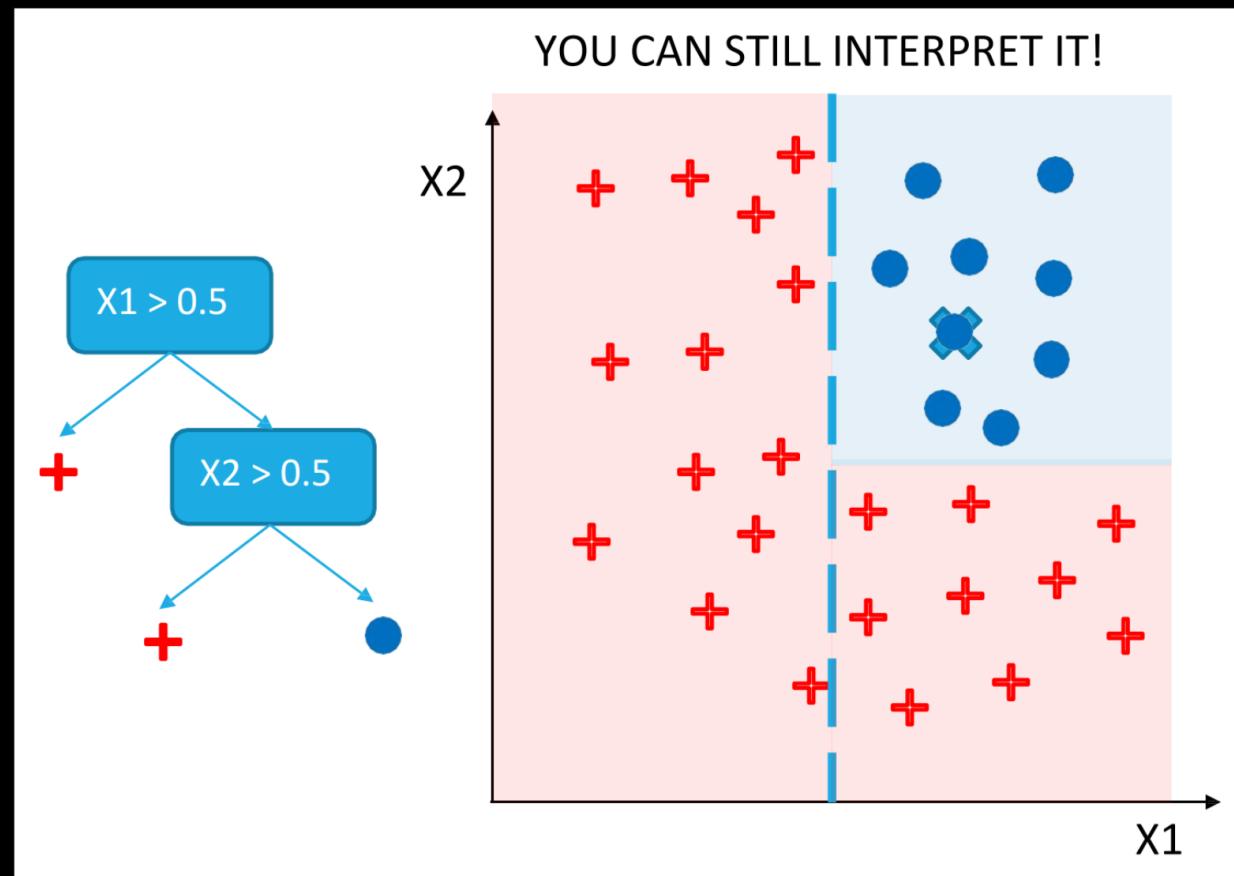
HISTORICAL DATA



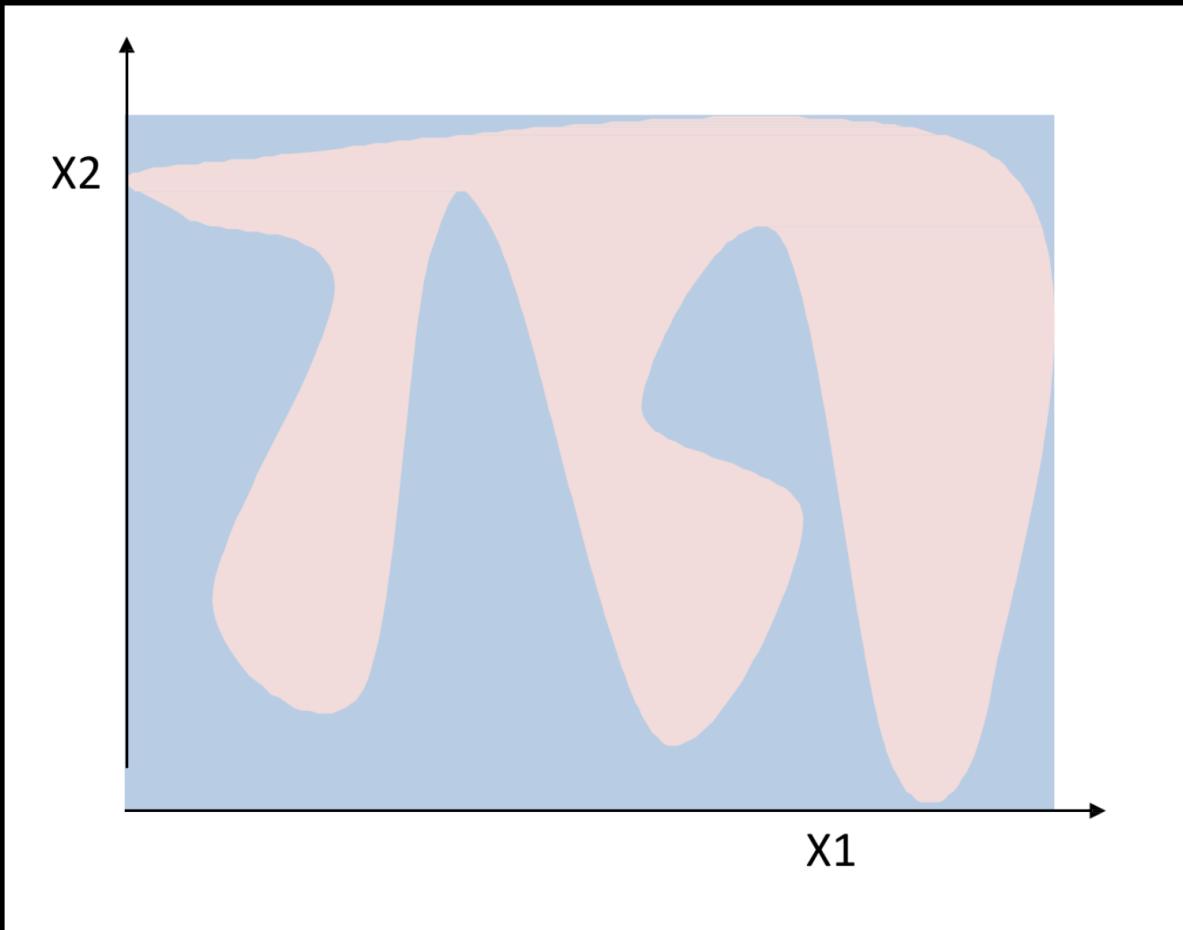
LINEAR CLASSIFIERS



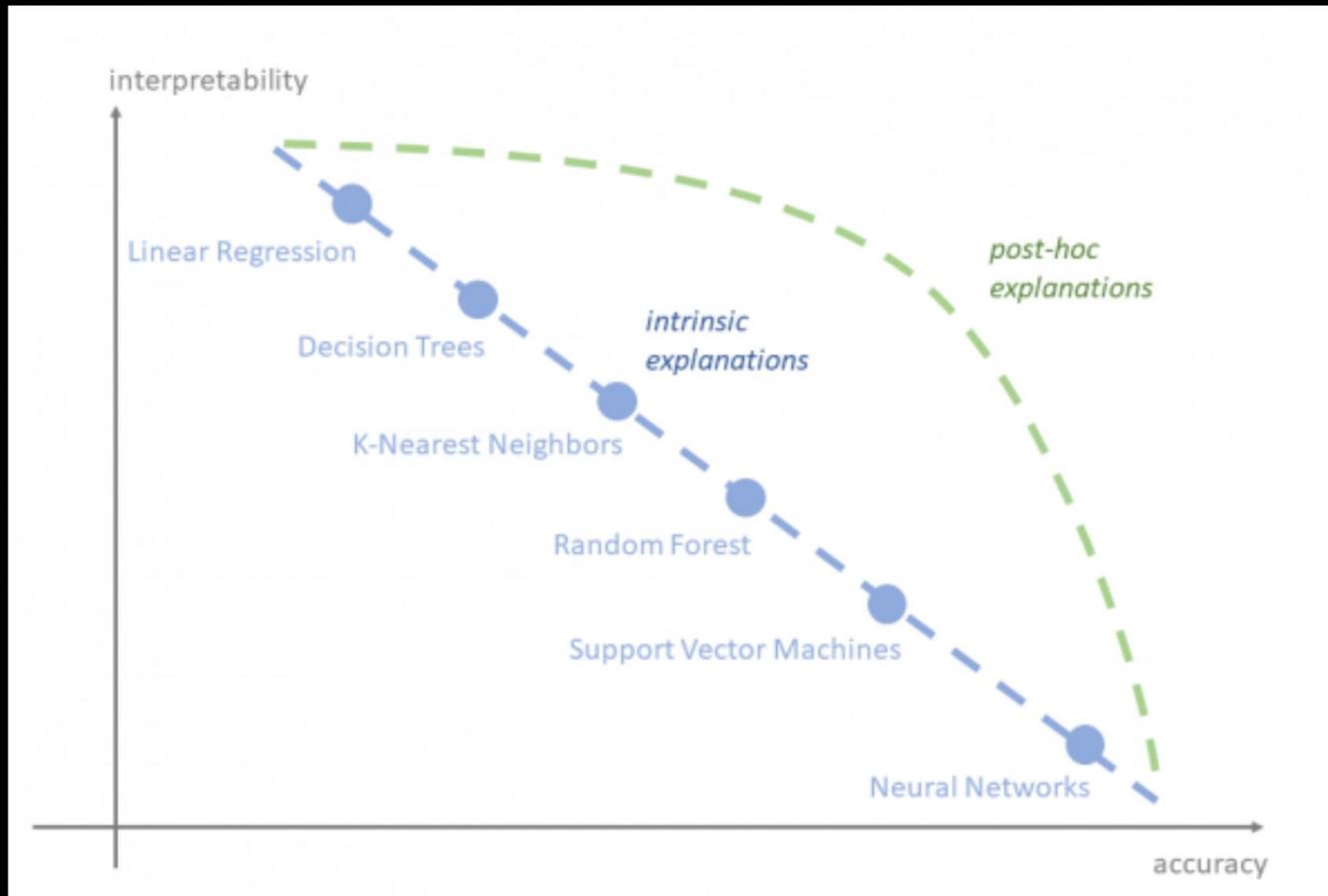
DECISION TREE



More Complexity & More Dimensions

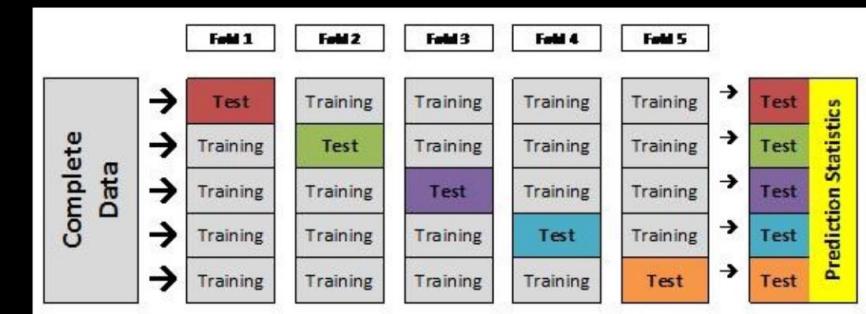
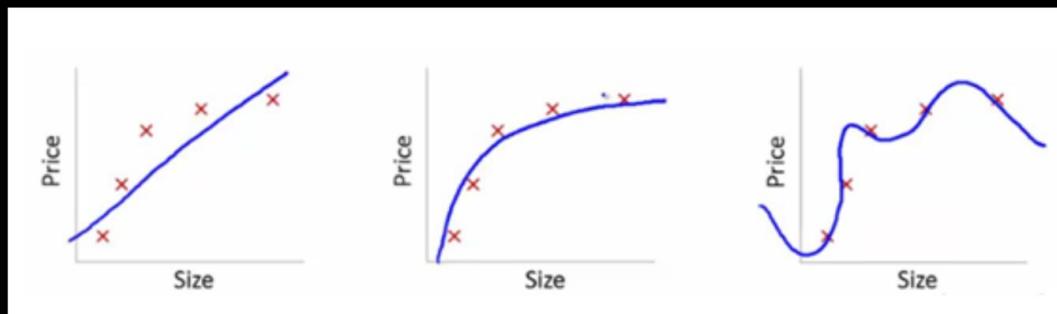


ACCURACY vs INTERPRETABILITY



WHAT PEOPLE DO TO GAIN TRUST?

- Restrict ourselves to only interpretable models
- Measure Accuracy (held-out data, cross validation)
- A/B testing



WHAT PEOPLE DO TO GAIN TRUST?



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME

Local Interpretable Model-agnostic Explanations



 Cornell University
Library

We gratefully acknowledge support from
the Simons Foundation
and member institutions

arXiv.org > cs > arXiv:1602.04938

Search or Article ID All fields

Computer Science > Machine Learning

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

Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin
(Submitted on 16 Feb 2016 (v1), last revised 9 Aug 2016 (this version, v3))

Despite widespread adoption, machine learning models remain mostly black boxes. Understanding the reasons behind predictions is, however, quite important in assessing trust, which is fundamental if one plans to take action based on a prediction, or when choosing whether to deploy a new model. Such understanding also provides insights into the model, which can be used to transform an untrustworthy model or prediction into a trustworthy one. In this work, we propose LIME, a novel explanation technique that explains the predictions of any classifier in an interpretable and faithful manner, by learning an interpretable model locally around the prediction. We also propose a method to explain models by presenting representative individual predictions and their explanations in a non-redundant way, framing the task as a submodular optimization problem. We demonstrate the flexibility of these methods by explaining different models for text (e.g. random forests) and image classification (e.g. neural networks). We show the utility of explanations via novel experiments, both simulated and with human subjects, on various scenarios that require trust: deciding if one should trust a prediction, choosing between models, improving an untrustworthy classifier, and identifying why a classifier should not be trusted.

Subjects: Machine Learning (cs.LG); Artificial Intelligence (cs.AI); Machine Learning (stat.ML)

Download:

- PDF
- Other formats

(license)

Current browse context:
cs.LG
< prev | next >
new | recent | 1602

Change to browse by:
cs
 cs.AI
stat
 stat.ML

References & Citations

- NASA ADS

DBLP - CS Bibliography

listing | bibtex
Marco Tulio Ribeiro
Sameer Singh
Carlos Guestrin

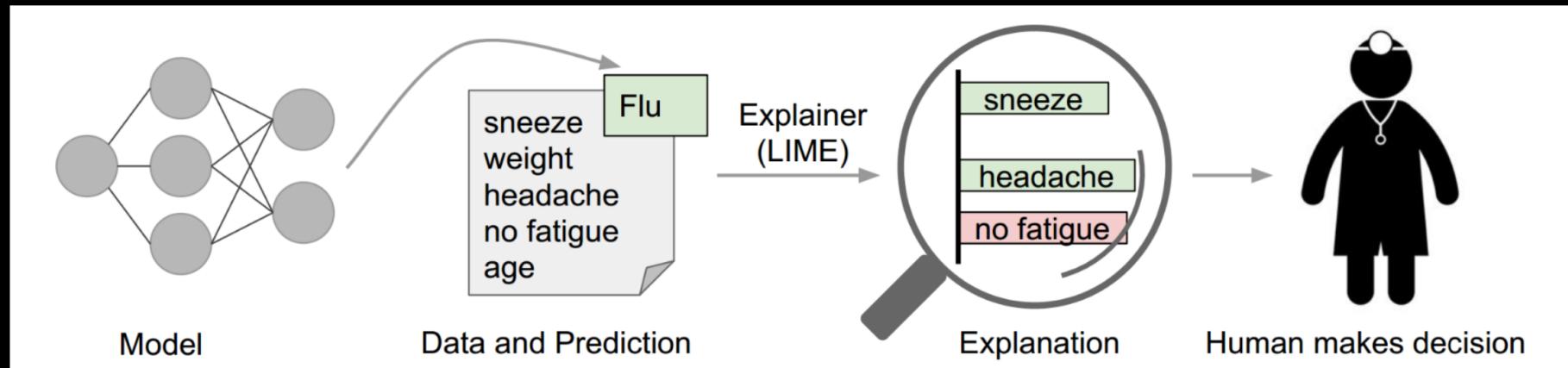
Google Scholar

Bookmark (what is this?)

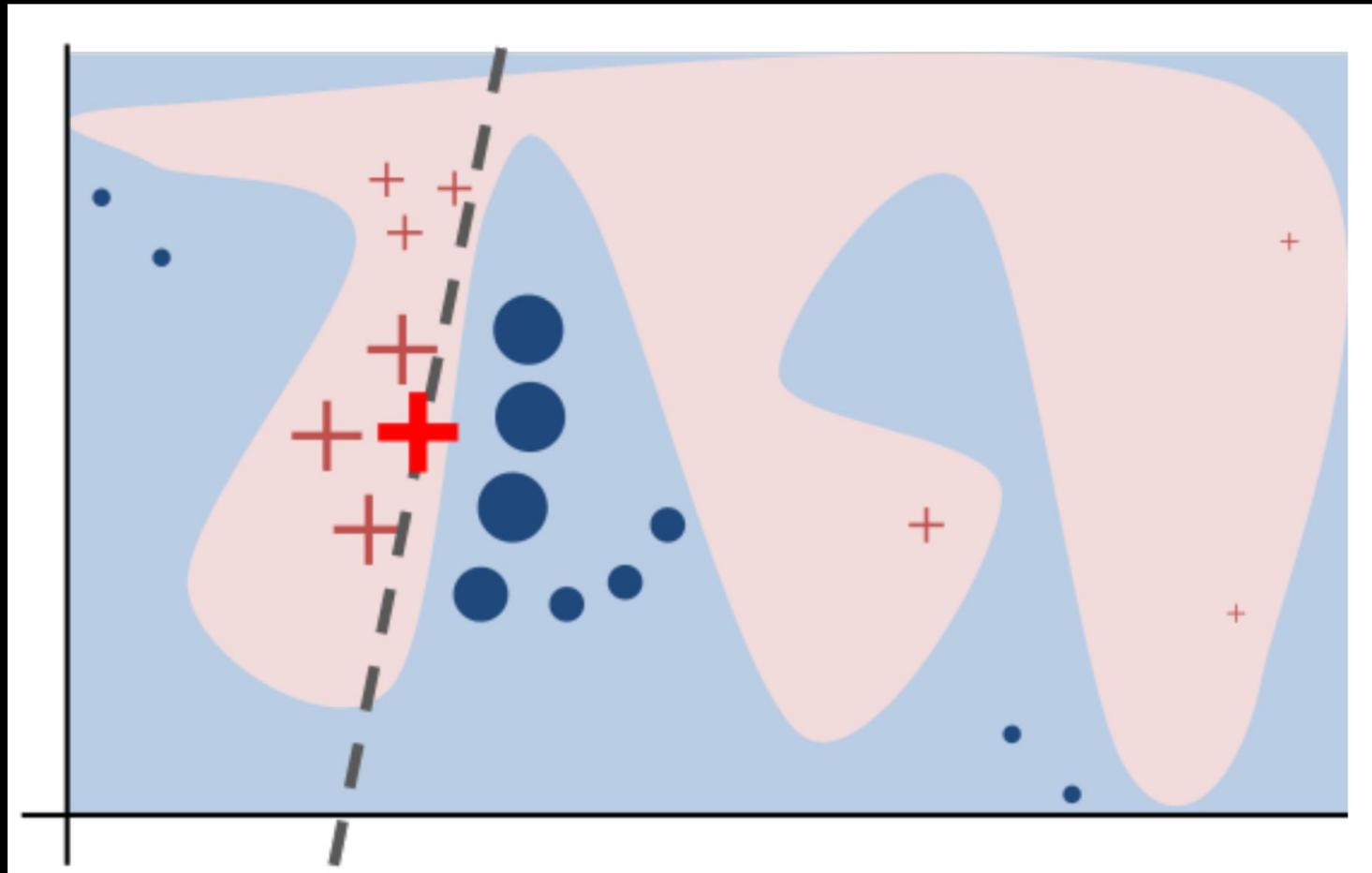


HOW IT
WORKS

Local Interpretable Model-agnostic Explanations



Local Interpretable Model-agnostic Explanations



Local Interpretable Model-agnostic Explanations



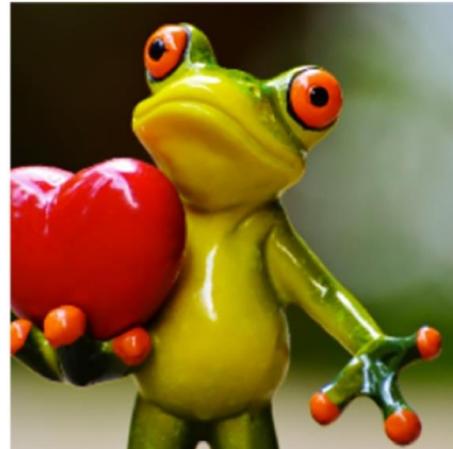
1. Permute data
2. Calculate distance between permutations and original observations
3. Make predictions on new data using complex model

Local Interpretable Model-agnostic Explanations



4. Pick m features best describing the complex model outcome from the data
5. Fit a simple model to the data with m features and similarity scores as weights
6. Feature weights from the simple model make explanations for the complex models local behavior

Local Interpretable Model-agnostic Explanations

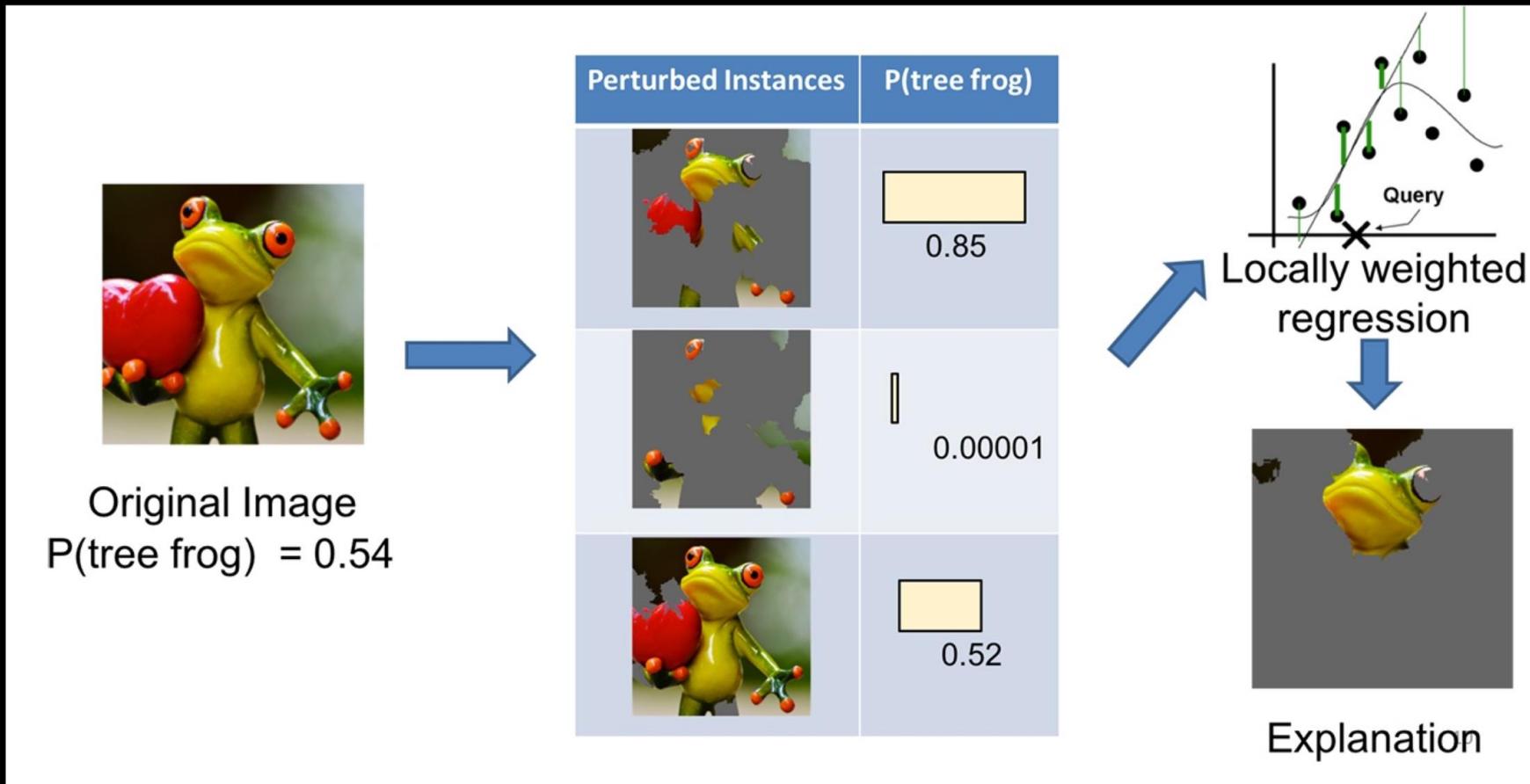


Original Image

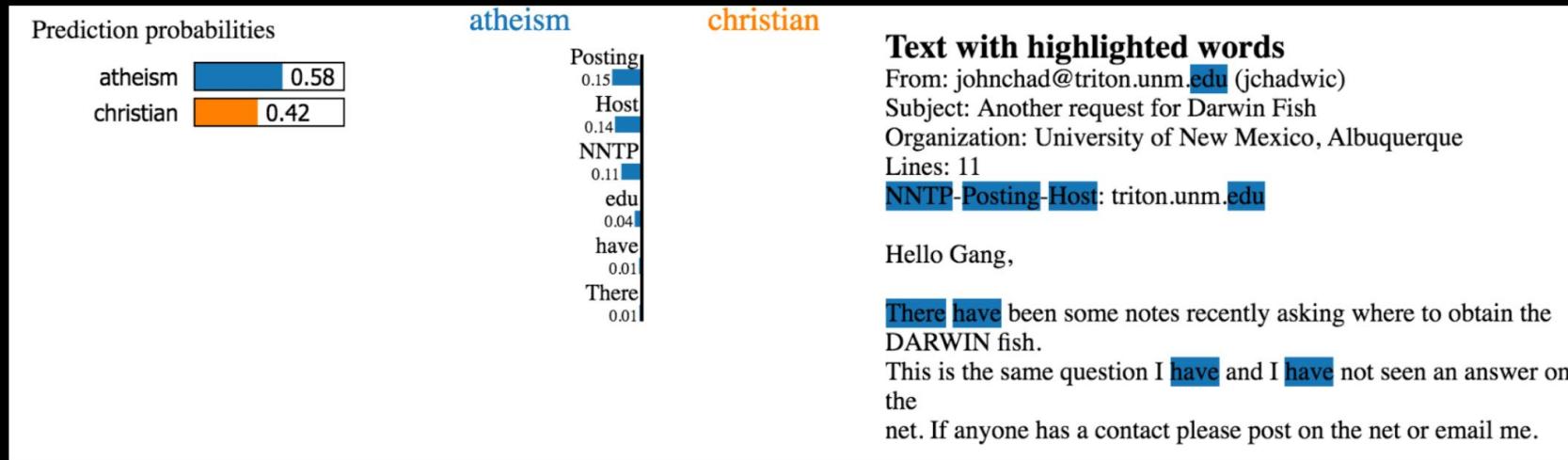


Interpretable
Components

Local Interpretable Model-agnostic Explanations



Local Interpretable Model-agnostic Explanations



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME

CREDIT SCORING



DATA

RANDOM FOREST

NEURAL NET

EXPLAIN



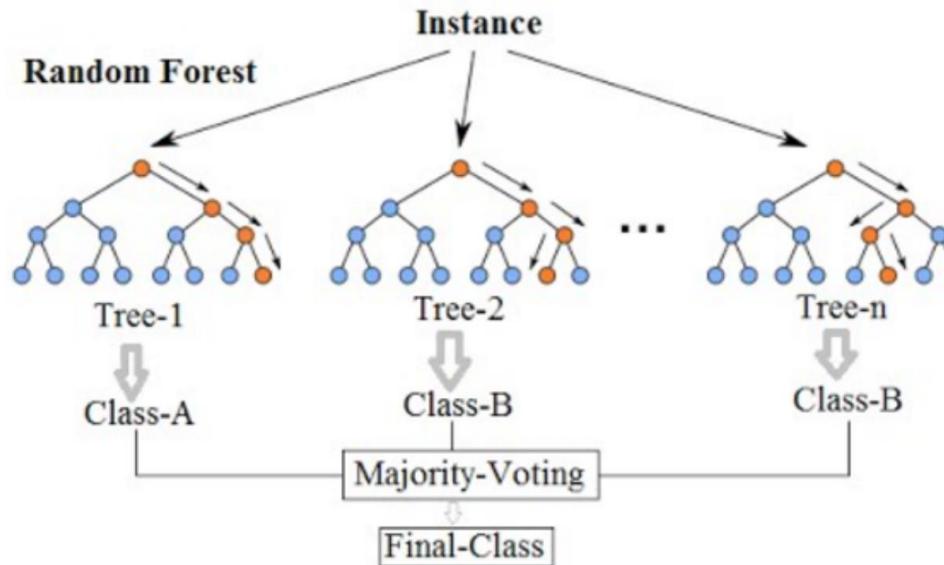
DATA

Field	Data	Description
Status	bad,good	Credit status
Home	ignore, other, owner, parents, priv, rent	Type of home ownership
Marital	divorced, married, separated, single, widow	Marital status
Records	no_rec, yes_rec	Existance of records
Job	fixed, freelance, others, partime	Type of job
Age		Client's age
Seniority		Job seniority (years)
Time		Maturity of requested loan (years)
Expenses		Expenses (in thousands \$)
Income		Income (in thousands \$)
Assets		Assets (in thousands \$)
Debt		Debt
Amount		Loan Amount (in thousands \$)
Price		Property price (in thousands \$)
FinRat		Financing Ratio
Savings		$\frac{\text{Income} - \text{Expenses} - (\text{Debt}/100)}{(\text{Amount} / \text{Time})}$

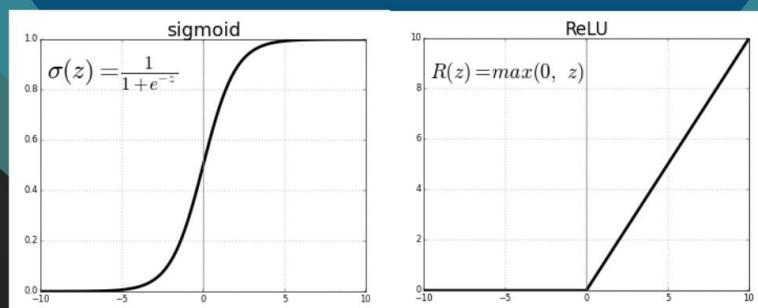
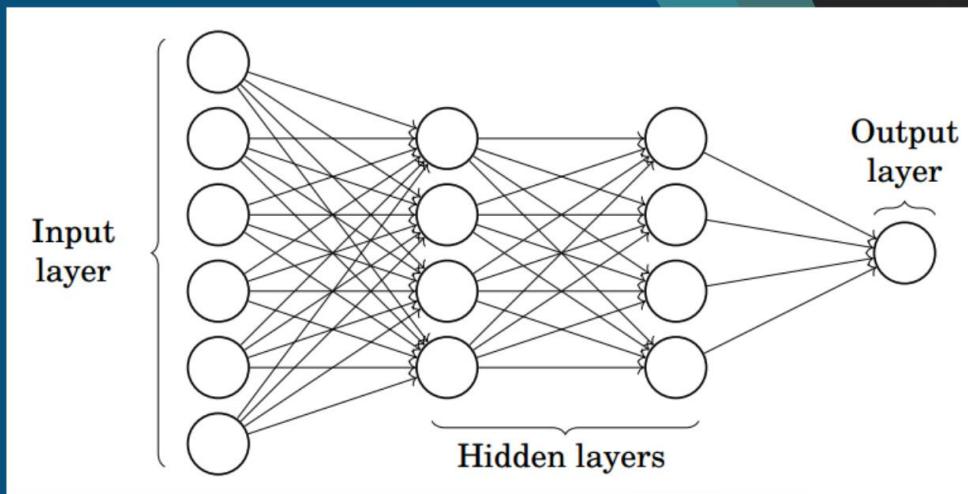
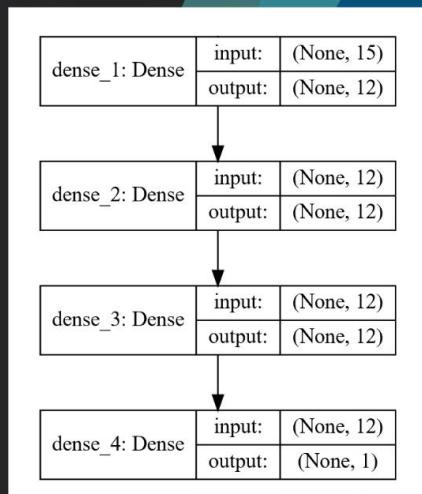


RANDOM FOREST

Random Forest Simplified

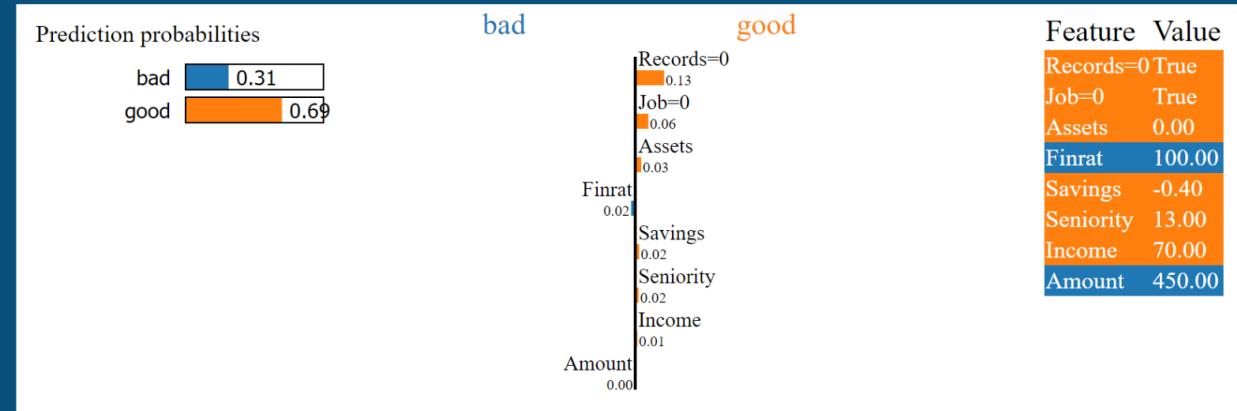


NEURAL NETWORK

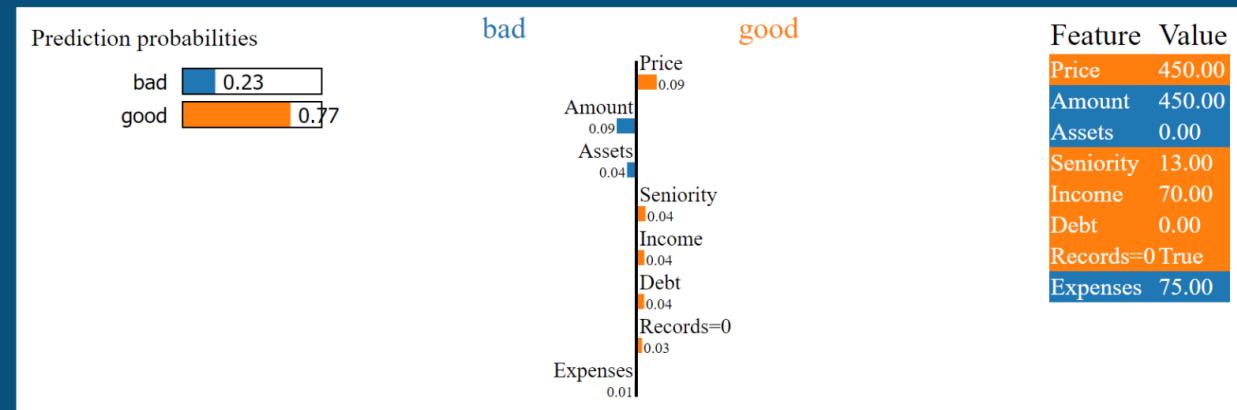


EXPLAIN - Random Forest vs Neural Net

Random
Forest



Neural Net
(MLP)



EXPLAIN - Compare Predictions

Status	Seniority	Home	Time	Age	Marital	Records	Job	Expenses	Income	Assets	Debt	Amount	Price	Finrat	Savings
bad	13	other	36	35	married	no_rec	fixed	75	70	0	0	450	450	100	-0.4
bad	13	other	36	35	married	no_rec	fixed	75	70	0	0	1200	1200	100	-0.4



MODEL RISK in AI/ML

WEAPONS OF MATH DESTRUCTION



ABOUT
ME

INTRO

BLACK BOX

FAST
BACKWARD

WRAP-UP

CREDIT
SCORING

LIME

WRAP-UP

TRUST
PREDICT
IMPROVE

ETHICS IN
AI

GDPR

OCBC

LET'S
CONNECT



TRUST, PREDICT, IMPROVE

Trust

How can we trust the predictions are correct?

Being able to interpret the explanations and compare classifiers based on them

Predict

How can we understand and predict the behavior?

Improved prediction of model behavior and time to make that assessment when explanations were provided

Improve

How do we improve it to prevent potential mistakes?

Non-ML experts with explanations vs ML experts without explanations



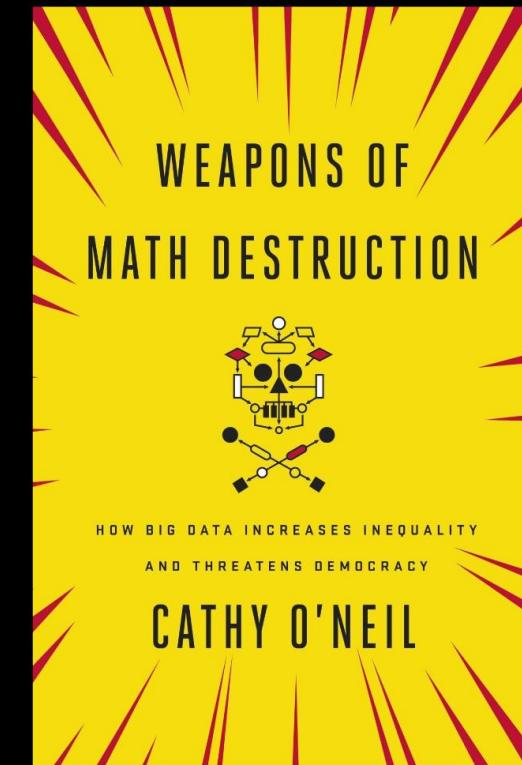
ETHICS IN AI

- Privacy
- Transparency
- Eliminate biases (gender, races, etc.)
- Explainability
- Accountability

“ THERE ARE ETHICAL CHOICES IN EVERY SINGLE ALGORITHM WE BUILD.



CATHY O'NEIL



GENERAL DATA PROTECTION REGULATION

Individual have now
the right to demand
an explanation of how
an AI system made a
decision that affects
them



KEEP
CALM
AND
COMPLY WITH
GDPR





Committed to being ethical and responsible

To both our customers and staff in our pursuit of being an AI enabled organisation

Ethical use of data & AI

OCBC is founding member of the MAS' committee to promote responsible and ethical use of AI

People & skillsets

OCBC's FutureSmart program trains and develops the digital skills of all employees of the OCBC Group

Committed to fair and transparent practices

Committed to developing our staff



LETS CONNECT!



www.linkedin.com/in/bertrndlenezet
https://github.com/lampalork/model_risk_ai_oct2018

