# Responsible Data Science

## Transparency & Interpretability

Auditing black-box models

*April 5, 2022*

**Prof. Julia Stoyanovich**

Center for Data Science &
Computer Science and Engineering
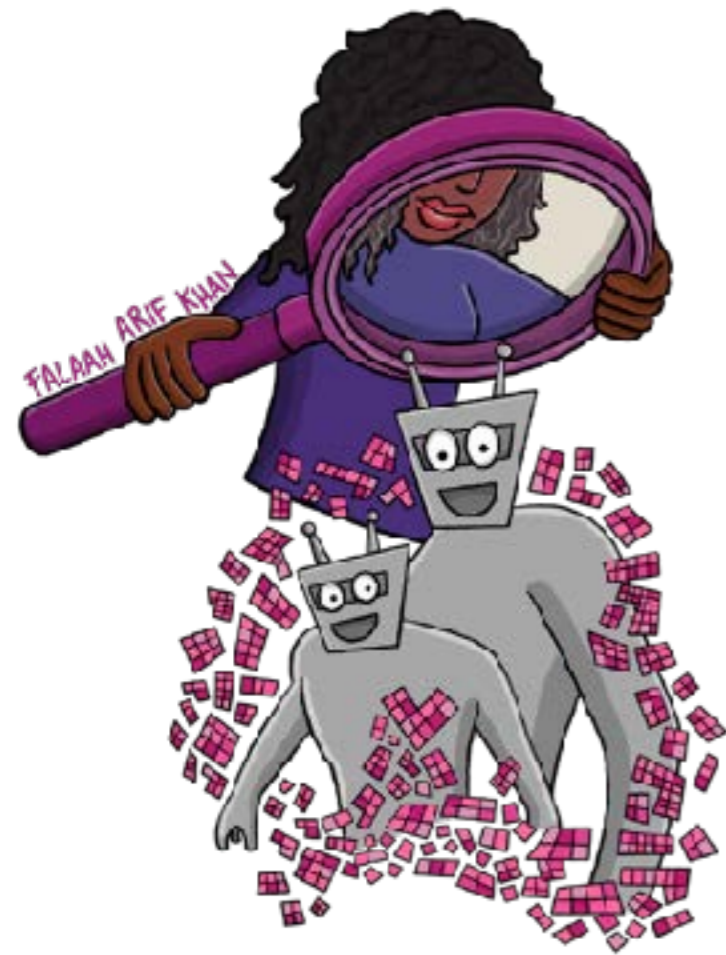New York University

NYU | TANDON SCHOOL OF ENGINEERING

NYU | Center for Data Science

r/ai

transparency, interpretability, explainability, intelligibility

responsible AI

agency, responsibility

# Interpretability for different stakeholders



**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

FALAAH ARIF KHAN

r/ai

# Staples discounts

**THE WALL STREET JOURNAL.**

WHAT THEY KNOW

## Websites Vary Prices, Deals Based on Users' Information

By Jennifer Valentino-DeVries, Jeremy Singer-Vine and Ashkan Soltani

December 24, 2012

**WHAT PRICE WOULD YOU SEE?**



It was the same Swingline stapler, on the same Staples.com website. But for Kim Wamble, the price was $15.79, while the price on Trude Frizzell's screen, just a few miles away, was $14.29.

A key difference: where Staples seemed to think they were located.

A Wall Street Journal investigation found that the Staples Inc. website displays different prices to people after estimating their locations. More than that, **Staples appeared to consider the person's distance from a rival brick-and-mortar store**, either OfficeMax Inc. or Office Depot Inc. If rival stores were within 20 miles or so, Staples.com usually showed a discounted price.

https://www.wsj.com/articles/SB10001424127887323777204578189391813881534

# Staples discounts

## THE WALL STREET JOURNAL.

WHAT THEY KNOW

### Websites Vary Prices, Deals Based on Users' Information

*By Jennifer Valentino-DeVries, Jeremy Singer-Vine and Ashkan Soltani*

December 24, 2012

**WHAT PRICE WOULD YOU SEE?**

It was the same Sw...
same Staples.com...
was $15.79, while...
a few miles away,...

A key difference:...
located.

A Wall Street Journal investigation found that the Sta... Inc. website displays different prices to people after estimating their locations. More than that, **Staples appeared to consider the person's distance from a rival brick-and-mortar store**, either OfficeMax Inc. or Office Depot Inc. If rival stores were within 20 miles or so, Staples.com usually showed a discounted price.

**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

https://www.wsj.com/articles/SB10001424127887323777204578189391813881534

r/ai

# Online job ads

**July 2015**

**Samuel Gibbs**

Wednesday 8 July 2015 11.29 BST

Automated testing and analysis of company's advertising system reveals male job seekers are shown far more adverts for high-paying executive jobs

ⓘ One experiment showed that Google displayed adverts for a career coaching service for executive jobs 1,852 times to the male group and only 318 times to the female group. Photograph: Alamy

## Women less likely to be shown ads for high-paid jobs on Google, study shows

The AdFisher tool simulated job seekers that did not differ in browsing behavior, preferences or demographic characteristics, except in gender.

One experiment showed that Google displayed ads for a career coaching service for "$200k+" executive jobs **1,852 times to the male group and only 318 times to the female group**. Another experiment, in July 2014, showed a similar trend but was not statistically significant.

https://www.theguardian.com/technology/2015/jul/08/women-less-likely-ads-high-paid-jobs-google-study

r/ai

**theguardian**

**July 2015**

Samuel Gibbs

Wednesday 8 July 2015 11.29 BST

Automated testing and analysis of company's advertising system reveals male job seekers are shown far more adverts for high-paying executive jobs

One experiment showed that Google displayed adverts for a career coaching service for executive jobs 1,852 times to the male group and only 318 times to the female group. Photograph: Alamy

# Women less likely to be shown ads for high-paid jobs on Google, study shows

The AdFisher tool simulated job seekers that did not differ in browsing behavior demographic

One experiment ads for a care executive job **and only 318** Another exp similar trend

**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

https://www.theguardian.com/technology/2015/jul/08/women-less-likely-ads-high-paid-jobs-google-study

r/ai

# Instant Checkmate



https://www.technologyreview.com/s/510646/racism-is-poisoning-online-ad-delivery-says-harvard-professor/

# Nutritional labels



**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

https://www.wsj.com/articles/why-the-labels-on-your-food-are-changing-or-

https://www.wsj.com/articles/imagine-a-nutrition-labelfor-

https://www.wsj.com/articles/hiring-job-candidates-ai-11632244313

# This week's reading

## Algorithmic Transparency via Quantitative Input Influence:
### Theory and Experiments with Learning Systems

Anupam Datta    Shayak Sen    Yair Zick
Carnegie Mellon University, Pittsburgh, USA
{danupam, shayaks, yairzick}@cmu.edu

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

Marco Tulio Ribeiro    Sameer Singh    Carlos Guestrin
University of Washington    University of Washington    University of Washington
Seattle, WA 98105, USA    Seattle, WA 98105, USA    Seattle, WA 98105, USA
marcotcr@cs.uw.edu    sameer@cs.uw.edu    guestrin@cs.uw.edu

## A Unified Approach to Interpreting Model Predictions

Scott M. Lundberg    Su-In Lee
Paul G. Allen School of Computer Science    Paul G. Allen School of Computer Science
University of Washington    Department of Genome Sciences
Seattle, WA 98105    University of Washington
slund1@cs.washington.edu    Seattle, WA 98105
    suinlee@cs.washington.edu

# What are we explaining?

How does a system work?

How **well** does a system work?

What does a system do?

Why was I ___ (mis-diagnosed / not offered a discount / denied credit) ?

Are a system's decisions discriminatory?

Are a system's decisions illegal?

# But isn't accuracy sufficient?



How is accuracy measured?  FPR / FNR / …

Accuracy for whom: over-all or in sub-populations?

Accuracy over which data?

There is never 100% accuracy.  Mistakes for what reason?

r/ai

# Facebook's real-name policy

Shane Creepingbear is a member of the Kiowa Tribe of Oklahoma

**October 13, 2014**

**Tweet**

Shane Creepingbear @Creepingbear · Oct 13, 2014
Hey yall  today I was kicked off of Facebook for having a fake name.
Happy Columbus Day great job #facebook #goodtiming #racist
#ColumbusDay

↺ 17

**TIME**

## Facebook Thinks Some Native American Names Are Inauthentic

BY JOSH SANBURN    FEBRUARY 14, 2015

**February 14, 2015**

If you're Native American, Facebook might think your name is fake.

The social network has a history of telling its users that the names they're attempting to use aren't real. Drag queens and overseas human rights activists, for example, have experienced error messages and problems logging in in the past.

The latest flap involves Native Americans, including Dana Lone Hill, who is Lakota. Lone Hill recently wrote in a blog post that Facebook told her her name was not "authentic" when she attempted to log in.

r/ai

# Explanations based on features

- **LIME** (Local Interpretable Model-Agnostic Explanations): to help users trust a prediction, explain individual predictions

- **SP-LIME**: to help users trust a model, select a set of representative instances for which to generate explanations



features in green ("sneeze", "headache") support the prediction ("Flu"), while features in red ("no fatigue") are evidence against the prediction

**what if patient id appears in green in the list? - an example of "data leakage"**

[Ribeiro, Singh & Guestrin, 2016]

# LIME: Local explanations of classifiers

Three must-haves for a good explanation

**Interpretable** • Humans can easily interpret reasoning



Definitely
not interpretable



Potentially
interpretable

[Ribeiro, Singh & Guestrin, 2016]

# Explanations based on features

Three must-haves for a good explanation

| Interpretable | • Humans can easily interpret reasoning |
|---|---|
| Faithful | • Describes how this model actually behaves |



Learned model

Not faithful to model

[Ribeiro, Singh & Guestrin, 2016]

# Explanations based on features

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 |

Can explain this mess ☺

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# Key idea: Interpretable representation

"The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classier."

- LIME relies on a distinction between **features** and **interpretable data representations**; examples:

  - In text classification features are word embeddings; an interpretable representation is a vector indicating the presence of absence of a word

  - In image classification features encoded in a tensor with three color channels per pixel; an interpretable representation is a binary vector indicating the presence or absence of a contiguous patch of similar pixels

- **To summarize**: we may have some $d$ features and $d'$ interpretable components; interpretable models will act over domain $\{0, 1\}^{d'}$ - denoting the presence of absence of each of d' interpretable components

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# Fidelity-interpretability trade-off

"The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classier."

$$f : \mathbb{R}^d \to \mathbb{R} \qquad g \in G, \, dom(g) = \{0,1\}^{d'} \qquad \Omega(g)$$

classifier **model being explained**

**explanation model**

some class of interpretable models

measure of complexity of explanation **g**

$f(x)$ denotes the probability that **x** belongs to some class

$\pi_x$ is a **proximity measure** relative to **x**

we make no assumptions about **f** to remain model-agnostic: draw samples weighted by $\pi_x$

**explanation**

measures how unfaithful is **g** to **f** in the locality around **x**

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

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# Fidelity-interpretability trade-off

"The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classier."

1. sample points around +

[Ribeiro, Singh & Guestrin, 2016]

# Fidelity-interpretability trade-off

"The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classier."

1. sample points around ╋

2. use complex model $f$ to assign class labels

[Ribeiro, Singh & Guestrin, 2016]

# Fidelity-interpretability trade-off

"The overall goal of LIME is to identify an **interpretable** model over the *interpretable representation* that is **locally faithful** to the classier."

1. sample points around +

2. use complex model **f** to assign class labels

3. weigh samples according to $\pi_x$

4. learn simple model **g** according to samples

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# Example: text classification with SVMs



**94% accuracy, yet we shouldn't trust this classifier!**

[Ribeiro, Singh & Guestrin, 2016]

## Explaining Google's Inception NN



probabilities of the top-3 classes
and the super-pixels predicting each

P( ) = 0.32

P( ) = 0.24

P( ) = 0.21







Electric guitar - incorrect but
reasonable, similar fretboard

Acoustic guitar

Labrador

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# When accuracy is not enough

Train a neural network to predict wolf v. husky



| | | | | | |
|---|---|---|---|---|---|
| Predicted: wolf<br>True: wolf | Predicted: husky<br>True: husky | Predicted: wolf<br>True: wolf | Predicted: wolf<br>True: husky | Predicted: husky<br>True: husky | Predicted: wolf<br>True: wolf |

Only 1 mistake!!!

Do you trust this model?
How does it distinguish between huskies and wolves?

[Ribeiro, Singh & Guestrin, 2016]

# When accuracy is not enough

Explanations for neural network prediction



We've built a great snow detector… ☹

[Ribeiro, Singh & Guestrin, 2016]

Why should I trust you?

Explaining the predictions of any classifier

Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin

Check out our paper, and open source project at
https://github.com/marcotcr/lime

https://www.youtube.com/watch?v=hUnRCxnydCc

[Ribeiro, Singh & Guestrin, 2016]

# Auditing black-box models



images by Anupam Datta

[Datta, Sen & Zick, 2016]

# QII: Quantitative Input Influence

Goal: determine how much influence an input, or a set of inputs, has on a **classification outcome** for an individual or a group

**Transparency queries / quantities of interest**

**Individual:** Which inputs have the most influence in my credit denial?

**Group:** Which inputs have the most influence on credit decisions for women?

**Disparity:** Which inputs influence men getting more positive outcomes than women?

[Datta, Sen & Zick, 2016]

r/ai

# Running example

Consider lending decisions by a bank, based on gender, age, education, and income. **Does gender influence lending decisions?**

- Observe that 20% of women receive the positive classification.

- To check whether gender impacts decisions, take the input dataset and replace the value of gender in each input profile by drawing it from the uniform distribution: set gender in 50% of the inputs to female and 50% to male.

- If we still observe that 20% of female profiles are positively classified **after the intervention** - we conclude that gender does not influence lending decisions.

- Do a similar test for other features, one at a time. This is known as **Unary QII**

r/ai

How much influence do individual features have a given classifier's decision about an individual?



| Age | 23 | DENIED |
| Workclass | Private | |
| Education | 11th | |
| Marital Status | Never married | |
| Occupation | Craft repair | |
| Relationship to household income | Child | |
| Race | Asian-Pac Island | |
| Gender | Male | |
| Capital gain | $14344 | |
| Capital loss | $0 | |
| Work hours per week | 40 | |
| Country | Vietnam | |

income

images by Anupam Datta

[Datta, Sen & Zick, 2016]

r/ai

# Transparency report: Mr. Y

Explanations for superficially similar individuals can be different



| | |
|---|---|
| Age | 27 |
| Workclass | Private |
| Education | Preschool |
| Marital Status | Married |
| Occupation | Farming-Fishing |
| Relationship to household income | Other Relative |
| Race | White |
| Gender | Male |
| Capital gain | $41310 |
| Capital loss | $0 |
| Work hours per week | 24 |
| Country | Mexico |

DENIED

income

images by Anupam Datta

[Datta, Sen & Zick, 2016]

r/ai

# Unary QII

For a quantity of influence *Q* and an input feature *i*, the QII of *i* on *Q* is the difference in *Q* when *i* is changed via an **intervention**.



replace features with random values from the population, examine the distribution over outcomes

[Datta, Sen & Zick, 2016]

For a quantity of influence *Q* and an input feature *i*, the QII of *i* on *Q* is the difference in *Q* when *i* is changed via an **intervention**.



**intervening on one feature at a time will not have any effect**

images by Anupam Datta

[Datta, Sen & Zick, 2016]

# Marginal QII

- Not all features are equally important within a set.

- *Marginal QII*: Influence of age and income over only income.

$$\iota(\{age, income\}) - \iota(\{income\})$$

**Need to aggregate Marginal QII across all sets**

- But age is a part of many sets!

$$\iota(\{age\}) - \iota(\{\})$$

$$\iota(\{age, gender, job\}) - \iota(\{gender, job\})$$
$$\iota(\{age, gender\}) - \iota(\{gender\})$$

$$\iota(\{age, job\}) - \iota(\{job\})$$

$$\iota(\{age, gender, job\}) - \iota(\{gender, job\})$$

$$\iota(\{age, gender, income\}) - \iota(\{gender, income\})$$
$$\iota(\{age, gender, income, job\}) - \iota(\{gender, income, job\})$$

# Aggregating influence across sets

**Idea:** Use game theory methods: voting systems, revenue division

*"In voting systems with multiple agents with differing weights, voting power often does not directly correspond to the weights of the agents. For example, the US presidential election can roughly be modeled as a cooperative game where each state is an agent. The **weight of a state is the number of electors in that state** (i.e., the number of votes it brings to the presidential candidate who wins that state). Although states like California and Texas have higher weight, swing states like Pennsylvania and Ohio tend to have higher power in determining the outcome of elections."*

This paper uses the **Shapley value** as the aggregation mechanism

$$\varphi_i(N,v) = \mathbb{E}_\sigma[m_i(\sigma)] = \frac{1}{n!} \sum_{\sigma \in \Pi(N)} m_i(\sigma)$$

[Datta, Sen & Zick, 2016]

# Aggregating influence across sets

**Idea:** Use game theory methods: voting systems, revenue division

This paper uses the **Shapley value** as the aggregation mechanism

$$\varphi_i(N,v) = \mathbb{E}_\sigma[m_i(\sigma)] = \frac{1}{n!} \sum_{\sigma \in \Pi(N)} m_i(\sigma)$$

$v : 2^N \to \mathbb{R}$    influence of a set of features **S** on the outcome

$\varphi_i(N,v)$       influence of feature **i**, given the set of features **N = {1,…, n}**

$\sigma \in \Pi(N)$       a permutation over the features in set **N**

$m_i(\sigma)$       payoff corresponding to this permutation

[Datta, Sen & Zick, 2016]

r/ai

# QII summary

- A principled (and beautiful!) framework for determining the influence of a feature, or a set of features, on a decision

- Works for black-box models, with the assumption that the full set of inputs is available

- Accounts for correlations between features

- "Parametrizes" on what quantity we want to set (QII), how we intervene, how we aggregate the influence of a feature across sets

- Experiments in the paper: interesting results

- Also in the paper: a discussion of **transparency under differential privacy**

r/ai

# SHAP: Shapley Additive Explanations

A unifying framework for interpreting predictions with "additive feature attribution methods", including LIME and QII, for **local explanations**



https://www.youtube.com/watch?v=wjd1G5bu_TY

[Lundberg & Lee, 2017]

# SHAP: Shapley Additive Explanations

A unifying framework for interpreting predictions with "**additive feature attribution methods**", including LIME and QII, for **local explanations**

- The best explanation of a **simple model** is the model itself: the explanation is both accurate and interpretable. For complex models we must use a simpler **explanation model** — an interpretable approximation of the original model.

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

**model being explained**

$$g \in G, \, dom(g) = \{0,1\}^{d'}$$

**explanation model** from a class of interpretable models, over a set of **simplified features**

- **Additive feature attribution methods** have an explanation model that is a linear function of binary variables

[Lundberg & Lee, 2017]

r/ai

# Additive feature attribution methods

**Additive feature attribution methods** have an explanation model that is a linear function of binary variables (simplified features)
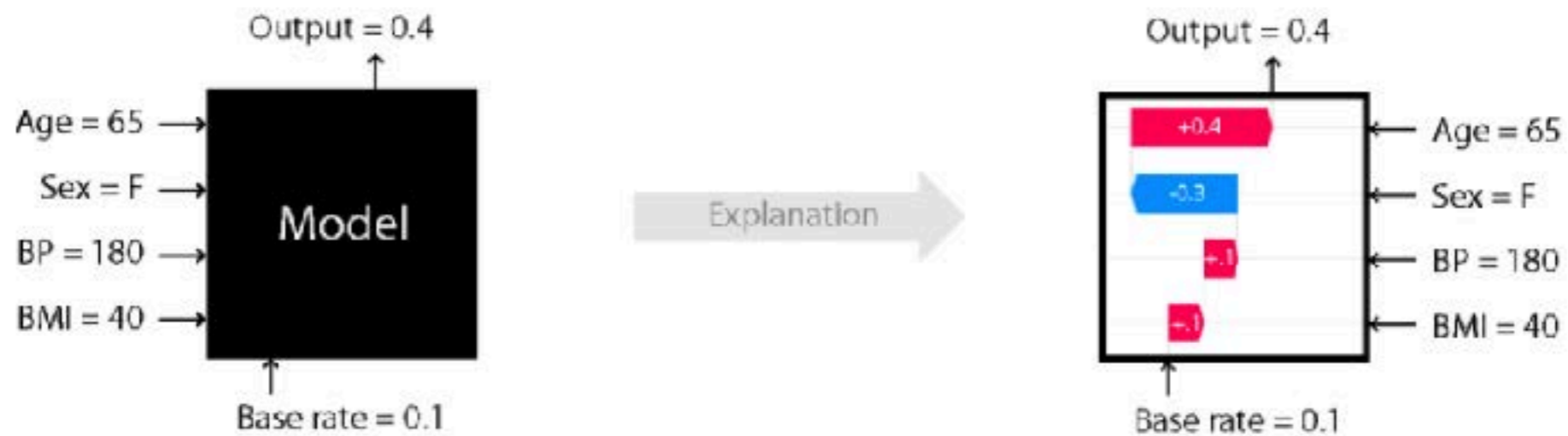
$$g(x') = \phi_0 + \sum_{i=1}^{d'} \phi_i x'_i \quad \text{where } x' \in \{0,1\}^{d'}, \text{and } \phi_i \in \mathbb{R}$$

Three properties guarantee a single unique solution — a unique allocation of Shapley values to each feature

1. **Local accuracy**: *g(x')* matches the original model *f(x)* when *x'* is the **simplified input** corresponding to x.

2. **Missingness**: if *x'$_i$* — the i[th] feature of simplified input *x'*— is missing, then it has no attributable impact for *x*

3. **Consistency** (**monotonicity**): if toggling off feature *i* makes a bigger (or the same) difference in model *f'(x)* than in model *f(x)*, then the weight (attribution) of *i* should be no lower in *f'(x)* than in *f(x)*

[Lundberg & Lee, 2017]

r/ai

# Additive feature attribution methods



https://github.com/slundberg/shap

[Lundberg & Lee, 2017]