# Transparency & Interpretability

**Prof. Julia Stoyanovich**
New York University

stoyanovich@nyu.edu
@stoyanoj

stoyanovich@nyu.edu
@stoyanoj

NYU | TANDON SCHOOL OF ENGINEERING

NYU | Center for Data Science

r/ai

# Supplementary reading

DOI:10.1038/s42256-020-0171-8

DOI:10.1145/3282486

responsible AI

transparency, interpretability, explainability, intelligibility

agency, responsibility

r/ai

# Interpretability for different stakeholders

**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

FALAAH ARIF KHAN

r/ai

# ADS in medical imaging

**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

FACEBOOK AI    NYU Langone Health

## fastMRI
Accelerating MR Imaging with AI

## What is fastMRI?

fastMRI is a collaborative research project between Facebook AI Research (FAIR) and NYU Langone Health. The aim is to investigate the use of AI to make MRI scans up to 10 times faster.

By producing accurate images from under-sampled data, AI image reconstruction has the potential to improve the patient's experience and to make MRIs accessible for more people.

To enable the broader research community to participate in this important project, NYU Langone Health has released fully anonymized raw data and image datasets. Visit our github repository, which contains baseline reconstruction models and PyTorch data loaders for the fastMRI dataset.

r/ai

# ADS in hiring



**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?

**ACCOUNTANT**
**Acme Partners**

| Qualifications: | BS in accounting, GPA >3.0, Knowledge of financial and accounting systems and applications |
| --- | --- |
| Personal data to be analyzed: | An AI program could be used to review and analyze the applicant's personal data online, including LinkedIn profile, social media accounts and credit score. |
| Additional assessment: | AI-assisted personality scoring |

**ALERT:** Applicants for this position DO NOT have the option to selectively decline use of AI analysis for any of their personal data or to review and challenge the results of such analysis.

r/ai

# Nutritional labels for ADS?

## Ranking Facts

### Ingredients →

| Attribute | Importance | |
|---|---|---|
| PubCount | 1.0 | 🌡️ |
| CSRankingAllArea | 0.24 | 🌡️ |
| Faculty | 0.12 | 🌡️ |

Importance of an attribute in a ranking is quantified by the correlation coefficient between attribute values and items scores, computed by a linear regression model. Importance is high if the absolute value of the correlation coefficient is over 0.75, medium if this value falls between 0.25 and 0.75, and low otherwise.

### Diversity overall ❓

**DeptSizeBin** ≡          **Regional Code** ≡

● Large ● Small          ● NE ● W ● MW ● SA ● SC

### Fairness ❓ →

| DeptSizeBin | FA*IR | | Pairwise | | Proportion | |
|---|---|---|---|---|---|---|
| Large | Fair | ✅ | Fair | ✅ | Fair | ✅ |
| Small | Unfair | ❌ | Unfair | ❌ | Unfair | ❌ |

A ranking is considered unfair when the p-value of the corresponding statistical test falls below 0.05.

### ← Stability

| Top-K | Stability |
|---|---|
| Top-10 | Stable |
| Overall | Stable |

**comprehensible**: short, simple, clear

**consultative**: provide actionable info

**comparable:** implying a standard

**computable:** incrementally constructed

[Yang, Stoyanovich, Asudeh, Howe, Jagadish, Miklau (2018)]
[Stoyanovich, Howe (2019)]

r/ai

# 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?

r/ai

# 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?

# Explanations based on features

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



Model      Data and Prediction      Explanation      Human makes decision

**what if patient id appears in green in the list?**

[Ribeiro, Singh & Guestrin, 2016]

1. sample points around +
2. use original model to assign class labels

**Key ideas**

interpretable features

interpretable models

locally faithful explanations

r/ai

# LIME: Locally Interpretable Model-Agnostic Explanations

1. sample points around **+**
2. use original model to assign class labels
3. weigh points according to distance from **+**
4. learn interpretable model according to samples



**Key ideas**

interpretable features

interpretable models

locally faithful explanations

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# When accuracy is not enough

Train a neural network to predict wolf v. husky



Predicted: wolf / True: wolf
Predicted: husky / True: husky
Predicted: wolf / True: wolf
Predicted: wolf / True: husky
Predicted: husky / True: husky
Predicted: wolf / True: wolf

Only 1 mistake!!!

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

**A snow detector!**

Explanations for neural network prediction



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 great snow detector... ☹

[Ribeiro, Singh & Guestrin, 2016]

r/ai

# When accuracy is not enough

## 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]

# Auditing black-box models

images by Anupam Datta

User data → Credit Classifier → Decisions

DENIED

User data → Credit Classifier → Decisions

[Datta, Sen & Zick, 2016]

r/ai

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

images by Anupam Datta



| | |
|---|---|
| Age | 23 |
| 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 |

DENIED

income

[Datta, Sen & Zick, 2016]

r/ai

Explanations for superficially similar individuals can be different

images by Anupam Datta



| | |
|---|---|
| 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

[Datta, Sen & Zick, 2016]
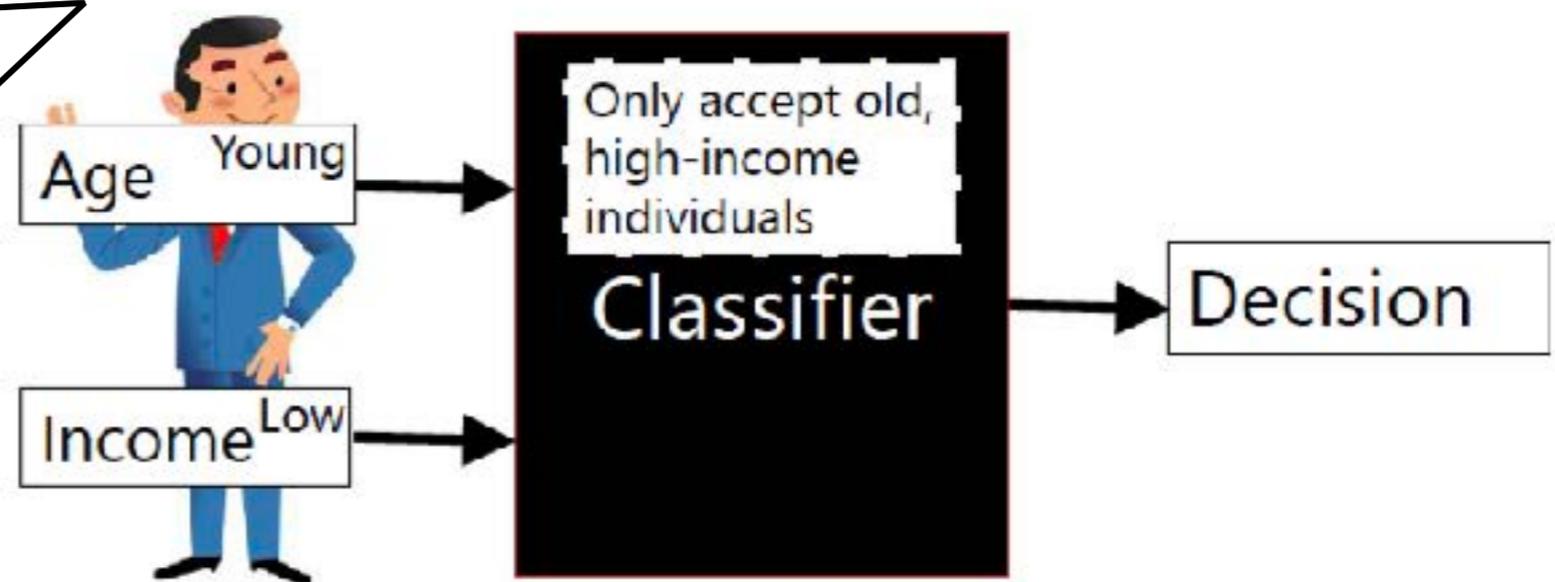
r/ai

# QII: Quantitative Input Influence

images by Anupam Datta

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]

r/ai

# QII: Quantitative Input Influence

images by Anupam Datta

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**.

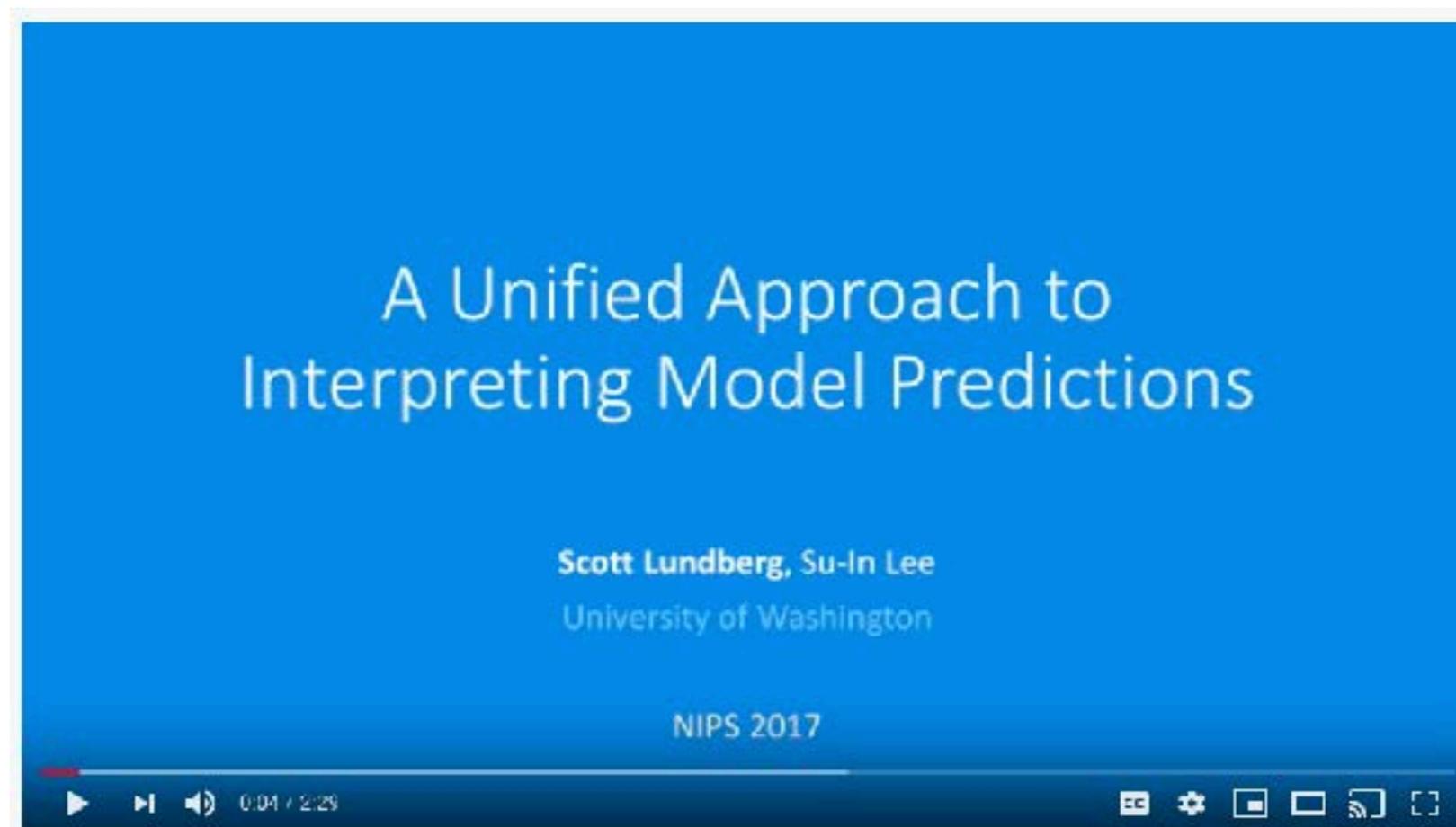**Key ideas**

**intervene** on an input feature, measure its **importance**

aggregate feature importance using its **Shapley value**



in this case, intervening on one feature at a time will have no effect

[Datta, Sen & Zick, 2016]

# SHAP: Shapley Additive Explanations

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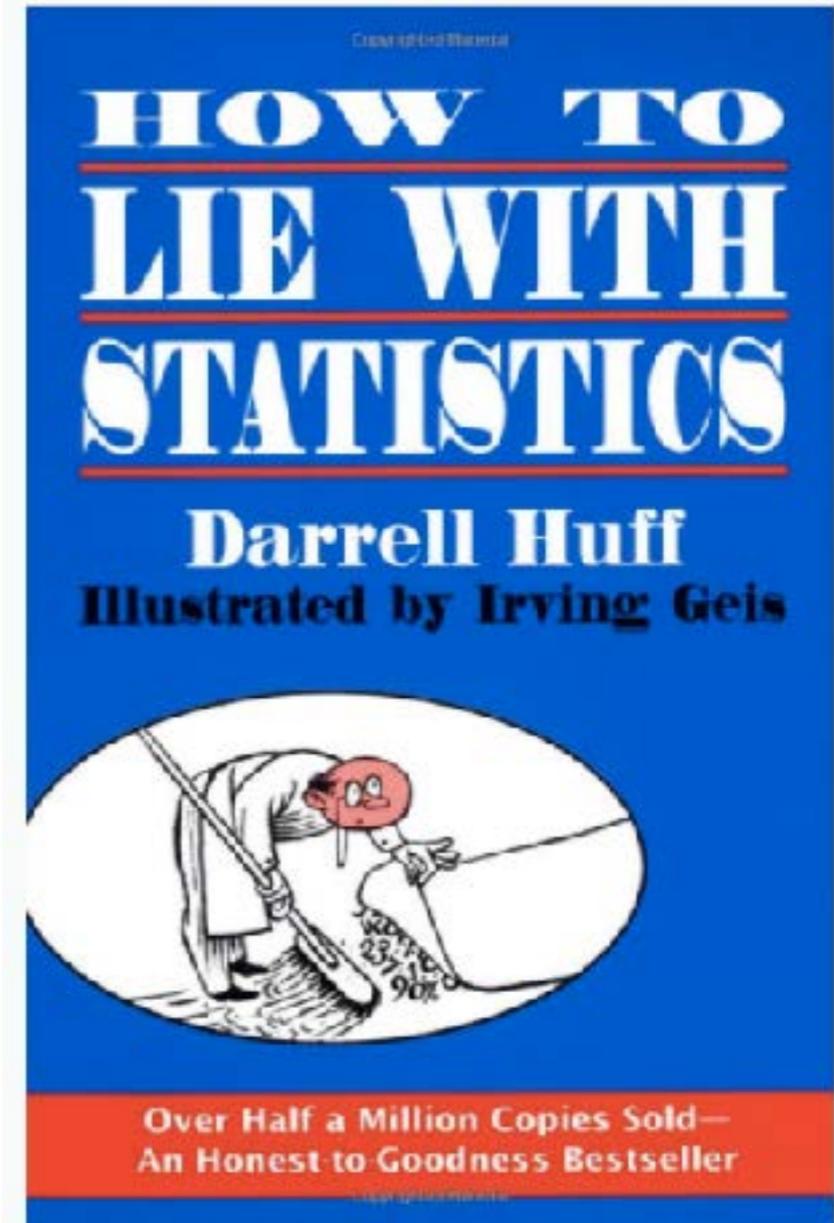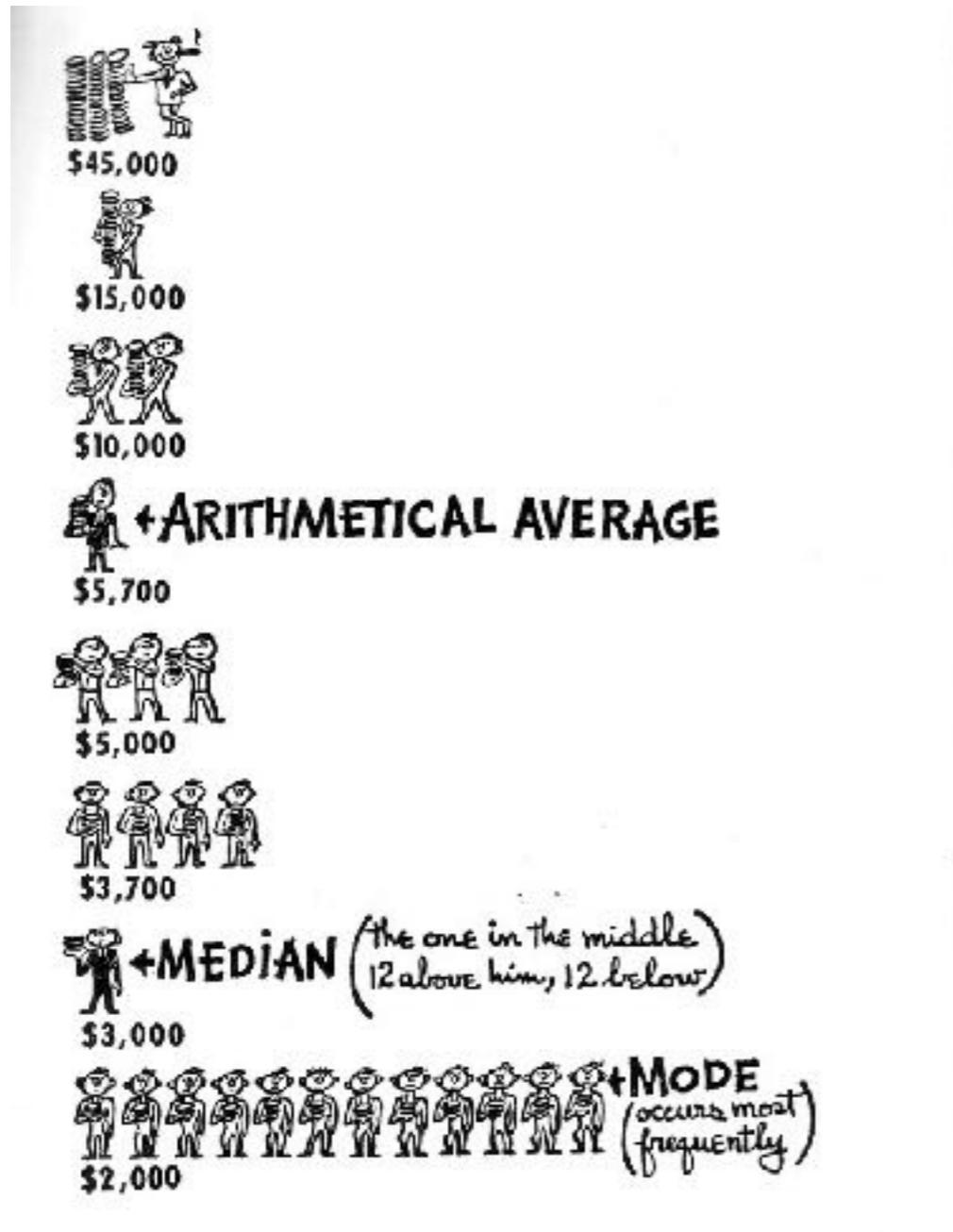


https://www.youtube.com/watch?v=wjd1G5bu_TY

[Lundberg & Lee, 2017]
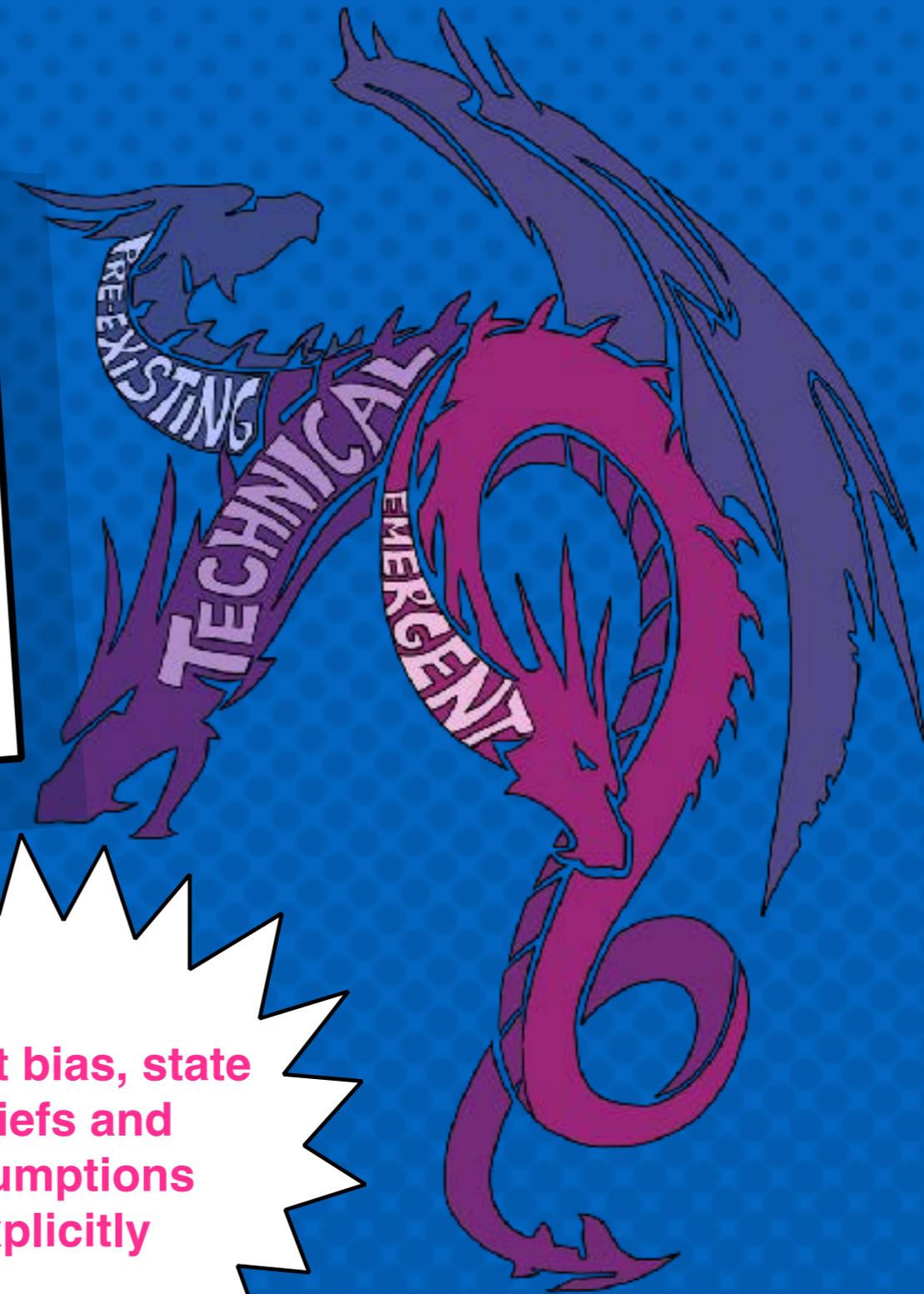
# Explaining the data

# The well-chosen average

# Explaining bias

**Pre-existing** is independent of an algorithm and has origins in society

**Technical** is introduced or exacerbated by the technical properties of an ADS

**Emergent** arises due to context of use

to fight bias, state beliefs and assumptions explicitly

[Friedman & Nissenbaum (1996)]

r/ai

# Explaining the models

**What** are we explaining?

To **Whom** are we explaining?

**Why** are we explaining?



FALAAH ARIF KHAN

r/ai

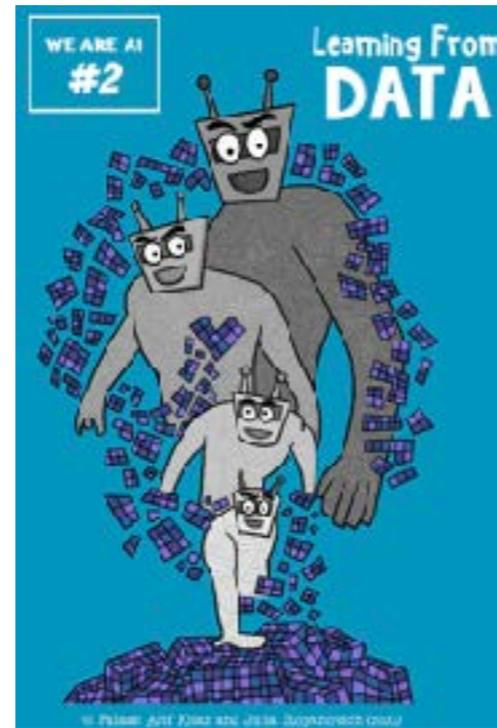dataresponsibly.github.io/we-are-ai

# AI comics for the general public



dataresponsibly.github.io/we-are-ai/comics

# Scientific comics on AI