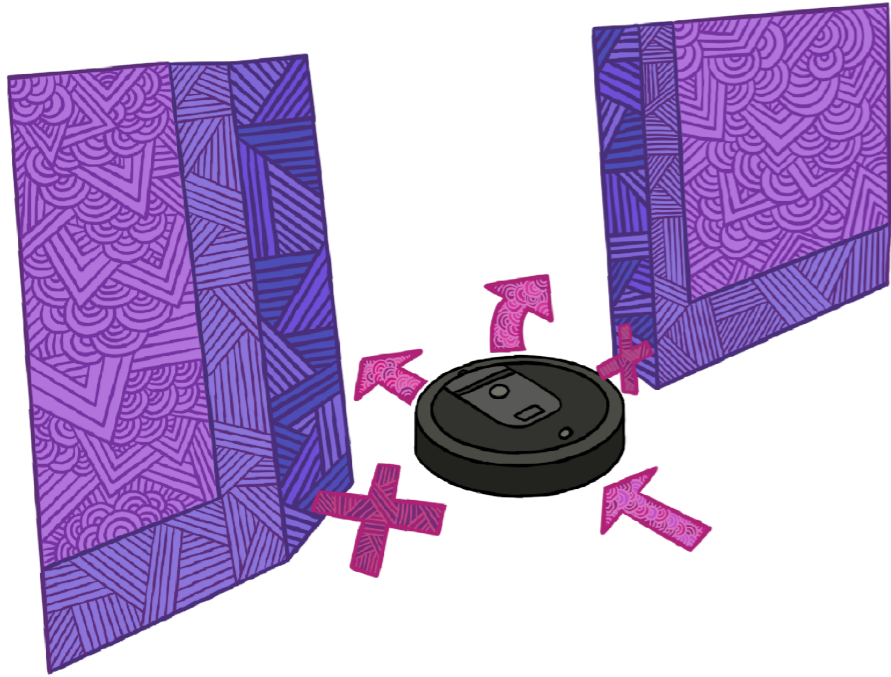# Responsible Data Management

**Julia Stoyanovich**

Computer Science and Engineering
Center for Data Science
Center for Responsible AI
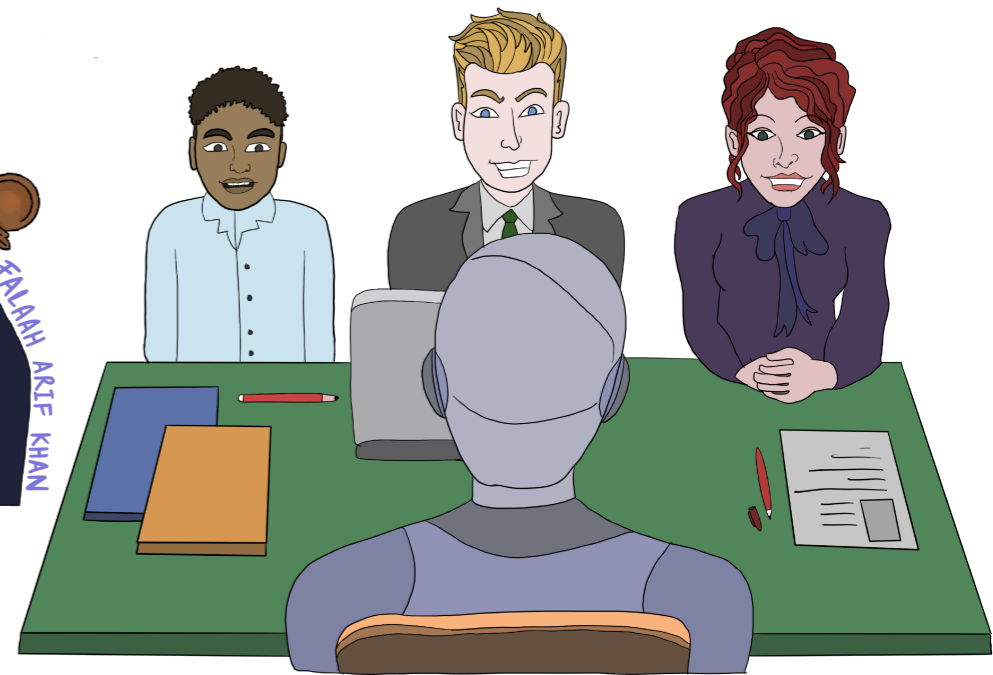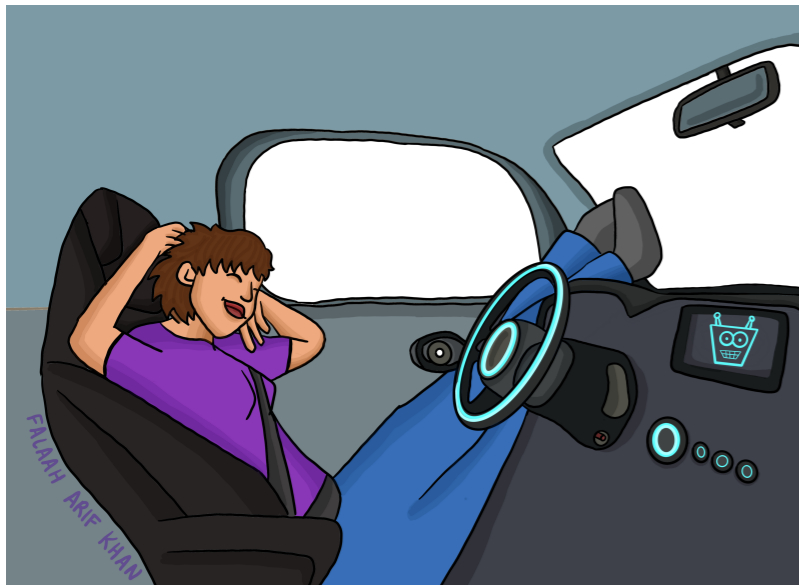Visualization and Data Analytics Center
New York University

# Machines make mistakes

FALAAH ARIF KHAN

# Harms can be cumulative



FALAAH ARIF KHAN

Sourcing

Screening

Interviewing

Background Check

Offers

# Racial bias in resume screening

## Are Emily and Greg More Employable Than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination

**September 2004**

Marianne Bertrand

Sendhil Mullainathan

**We study race in the labor market by sending fictitious resumes to help-wanted ads in Boston and Chicago newspapers.** To manipulate perceived race, resumes are randomly assigned African-American- or White-sounding names. **White names receive 50 percent more callbacks for interviews.** Callbacks are also more responsive to resume quality for White names than for African-American ones. The racial gap is uniform across occupation, industry, and employer size. We also find little evidence that employers are inferring social class from the names. Differential treatment by race still appears to still be prominent in the U. S. labor market.

r/ai center for responsible ai

# Bias in algorithmic hiring

**theguardian** July 2015

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

**REUTERS** October 2018

Amazon scraps secret AI recruiting tool that showed bias against women

THE WALL STREET JOURNAL. September 2014

Are Workplace Personality Tests Fair?

Growing Use of Tests Sparks Scrutiny Amid Questions of Effectiveness and Workplace Discrimination

The New York Times March 2021

## We Need Laws to Take On Racism and Sexism in Hiring Technology

Artificial intelligence used to evaluate job candidates must not become a tool that exacerbates discrimination.

**MIT Technology Review** February 2013

## Racism is Poisoning Online Ad Delivery, Says Harvard Professor

r/ai center for responsible ai

THE NEW YORK CITY COUNCIL
Corey Johnson, Speaker

**December 11, 2021**

This law requires that a **bias audit** be conducted on an automated employment decision tool prior to the use of said tool. The bill also requires that candidates or employees **be notified about the use of such tools** in the assessment or evaluation for hire or promotion before these tools are used, as well as **be notified about the job qualifications and characteristics that will be used** by the tool. Violations of the provisions of the bill are subject to a civil penalty.

r/ai center for responsible ai

# Personality prediction in hiring

# Algorithmic personality tests

**Input**: resume or LinkedIn handle (both systems) or Twitter (Humantic AI)

**Output**: a personality profile + a job fit score (Crystal) or match score (Humantic AI)

# Stability audit framework

Check for updates

## An external stability audit framework to test the validity of personality prediction in AI hiring

Alene K. Rhea[1,2] · Kelsey Markey[1,2] · Lauren D'Arinzo[1,2,3] ·
Hilke Schellmann[4] · Mona Sloane[2] · Paul Squires[5] · Falaah Arif Khan[1,2] ·
Julia Stoyanovich[1,2,6]

https://link.springer.com/article/10.1007/s10618-022-00861-0

# Stability audit framework

# Stability audit framework



| Facet | Crystal | Humantic |
|---|:---:|:---:|
| Resume file format | ✗ | ✓ |
| LinkedIn URL in resume | ? | ✗ |
| Source context | ✗ | ✗ |
| Algorithm-time / immediate | ✓ | ✓ |
| Algorithm-time / 31 days | ✓ | ✗ |
| Participant-time / LinkedIn | ✗ | ✗ |
| Participant-time / Twitter | N/A | ✓ |

**Pre-existing bias** has origins in society

FALAAH ARIF KHAN

**Pre-existing bias** has origins in society

**Pre-existing bias** has origins in society

Pre-existing bias has origins in society

FALAAH ARIF KHAN

r/ai center for responsible ai

# Diverse balanced ranking

## Goals

**diversity**: pick **k = 4** candidates, including 2 of each gender, and at least one per race

**utility**: maximize the total score of selected candidates

**score = 373**

**score = 372**

|        | Male     |          | Female   |          |
|--------|----------|----------|----------|----------|
| White  | A (99)   | B (98)   | C (96)   | D (95)   |
| Black  | E (91)   | F (91)   | G (90)   | H (89)   |
| Asian  | I (87)   | J (87)   | K (86)   | L (83)   |

## Problem

picked the best White and male candidates (A, B) but did not pick the best Black (E, F), Asian (I, J), or female (C, D) candidates

## Beliefs

scores are more informative within a group than across groups - **effort is relative to circumstance**

it is important to **reward effort**

[Yang, Gkatzelis, Stoyanovich (2019)]

# From beliefs to interventions

**Fairness for female candidates**    **83 / 95 = 0.91**

| C | D | G | H | K | L |
|---|---|---|---|---|---|
| 95 | 95 | 90 | 86 | 83 | 83 |

highest-scoring skipped

lowest-scoring selected

**BEFORE: diversity constraints only**



**AFTER: diversity and fairness constraints**



## Beliefs

scores are more informative within a group than across groups - **effort is relative to circumstance**

it is important to **reward effort**

[Yang, Gkatzelis, Stoyanovich (2019)]

r/ai center for responsible ai

# Fairness in Ranking, Part I: Score-based Ranking

MEIKE ZEHLIKE, Humboldt University of Berlin, Max Planck Institute for Software Systems, and Zalando Research, Germany

KE YANG, New York University, NY, and University of Massachusetts, Amherst, MA, USA

JULIA STOYANOVICH, New York University, NY, USA

In the past few years, there has been much work on incorporating fairness requirements into algorithmic rankers, with contributions coming from the data management, algorithms, information retrieval, and recommender systems communities. In this survey we give a systematic overview of this work, offering a broad perspective that connects formalizations and algorithmic approaches across subfields. An important contribution of our work is in developing a common narrative around the value frameworks that motivate specific fairness-enhancing interventions in ranking. This allows us to unify the presentation of mitigation objectives and of algorithmic techniques to help meet those objectives or identify trade-offs.

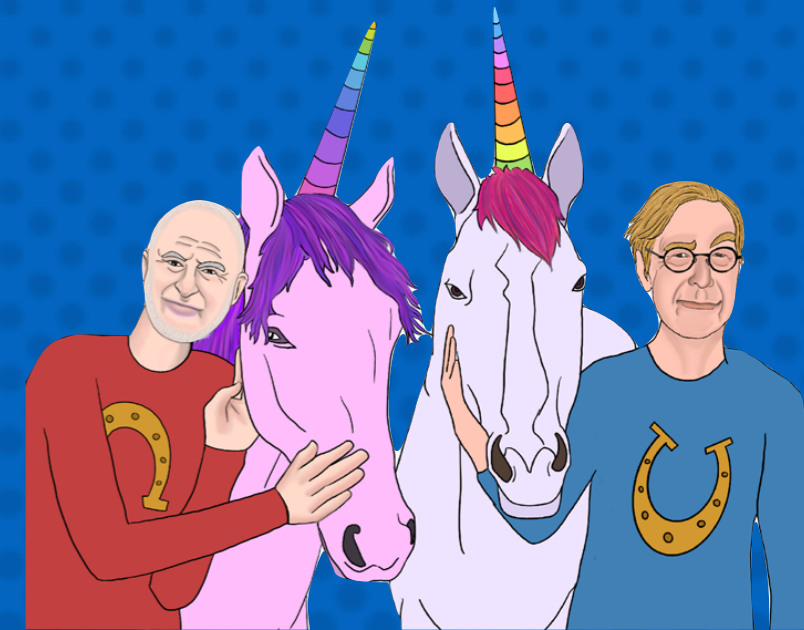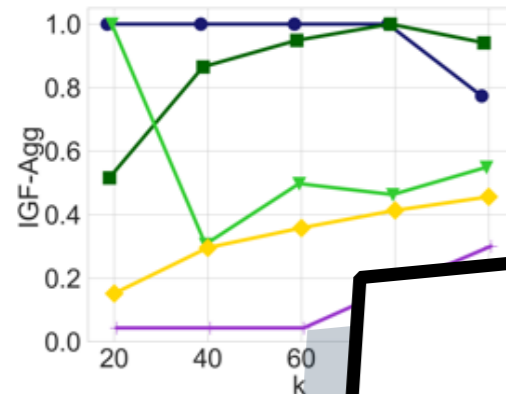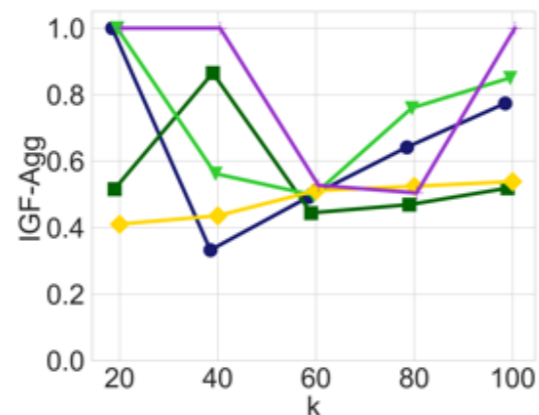In this first part of this survey, we describe four classification frameworks for fairness-enhancing interventions, along which we relate the technical methods surveyed in this paper, discuss evaluation datasets, and present technical work on fairness in score-based ranking. In the second part of this survey, we present methods that incorporate fairness in supervised learning, and also give representative examples of recent v
frameworks for fair score-based rankin
ranking methods.

# Fairness in Ranking, Part II: Learning-to-Rank and Recommender Systems

MEIKE ZEHLIKE, Humboldt University of Berlin, Max Planck Institute for Software Systems, and Zalando Research, Germany

KE YANG, New York University, NY, and University of Massachusetts, Amherst, MA, USA

JULIA STOYANOVICH, New York University, NY, USA

In the past few years, there has been much work on incorporating fairness requirements into algorithmic rankers, with contributions coming from the data management, algorithms, information retrieval, and recommender systems communities. In this survey we give a systematic overview of this work, offering a broad perspective that connects formalizations and algorithmic approaches across subfields. An important contribution of our work is in developing a common narrative around the value frameworks that motivate specific fairness-enhancing interventions in ranking. This allows us to unify the presentation of mitigation objectives and of algorithmic techniques to help meet those objectives or identify trade-offs.
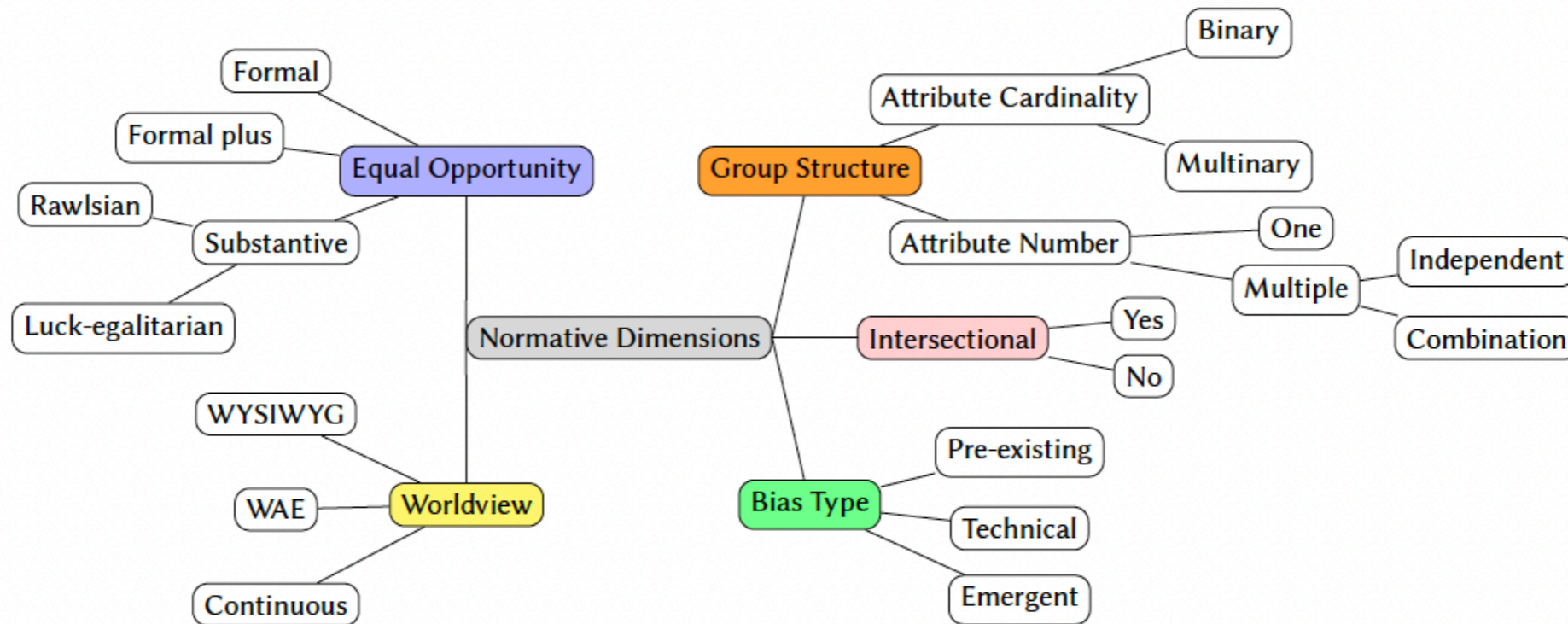
In the first part of this survey, we describe four classification frameworks for fairness-enhancing interventions, along which we relate the technical methods surveyed in this paper, discuss evaluation datasets, and present technical work on fairness in score-based ranking. In this second part of this survey, we present methods that incorporate fairness in supervised learning, and also give representative examples of recent work on fairness in recommendation and matchmaking systems. We also discuss evaluation frameworks for fair
air ranking methods.

[Zehlike, Yang, Stoyanovich (2022)]

**Technical bias** may be introduced or exacerbated by the technical properties of an ADS

# Model development lifecycle

**Goal**

design a model to predict an appropriate level of compensation for job applicants

**Problem**

accuracy is lower for middle-aged women - **a fairness concern!**

**now what?**

demographics

employment

split

interpolate missing

preprocess

select model

tune & validate

[Schelter, He, Khilnani, Stoyanovich (2020)]

center for responsible ai

# Missing value imputation

are values **missing at random** (e.g., *gender*, *age*, *years of experience*, *disability status* on job applications)?

are we ever interpolating **rare categories** (e.g., *Native American*)

are **all categories** represented (e.g., *non-binary gender*)?

"filtering" operations (like selection and join), **can arbitrarily change demographic group proportions**

select by zip code, country, years of C++ experience, others?

| age_group | county |
|-----------|--------|
| 60 | CountyA |
| 60 | CountyA |
| 20 | CountyA |
| 60 | CountyB |
| 20 | CountyB |
| 20 | CountyB |

→

| age_group | county |
|-----------|--------|
| 60 | CountyA |
| 60 | CountyA |
| 20 | CountyA |

66% vs 33%

50% vs 50%

r/ai center for responsible ai

# Data filtering

"filtering" operations (like selection and join), **can arbitrarily change demographic group proportions**

select by zip code, country, years of C++ experience, others?

# Data distribution debugging: mlinspect

**Potential issues in preprocessing pipeline:**

**1** Join might change proportions of groups in data

**2** Column 'age_group' projected out, but required for fairness

**3** Selection might change proportions of groups in data

**4** Imputation might change proportions of groups in data

**5** 'race' as a feature might be illegal!

**6** Embedding vectors may not be available for rare names!

**Python script for preprocessing, written exclusively with native pandas and sklearn constructs**

```python
# load input data sources, join to single table
patients = pandas.read_csv(…)
histories = pandas.read_csv(…)
data = pandas.merge([patients, histories], on=['ssn'])

# compute mean complications per age group, append as column
complications = data.groupby('age_group')
 .agg(mean_complications=('complications','mean'))
data = data.merge(complications, on=['age_group'])

# Target variable: people with frequent complications
data['label'] = data['complications'] >
    1.2 * data['mean_complications']

# Project data to subset of attributes, filter by counties
data = data[['smoker', 'last_name', 'county',
             'num_children', 'race', 'income', 'label']]
data = data[data['county'].isin(counties_of_interest)]

# Define a nested feature encoding pipeline for the data
impute_and_encode = sklearn.Pipeline([
  (sklearn.SimpleImputer(strategy='most_frequent')),
  (sklearn.OneHotEncoder())])
featurisation = sklearn.ColumnTransformer(transformers=[
  (impute_and_encode, ['smoker', 'county', 'race']),
  (Word2VecTransformer(), 'last_name')
  (sklearn.StandardScaler(), ['num_children', 'income']])

# Define the training pipeline for the model
neural_net = sklearn.KerasClassifier(build_fn=create_model())
pipeline = sklearn.Pipeline([
  ('features', featurisation),
  ('learning_algorithm', neural_net)])

# Train-test split, model training and evaluation
train_data, test_data = train_test_split(data)
model = pipeline.fit(train_data, train_data.label)
print(model.score(test_data, test_data.label))
```
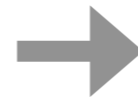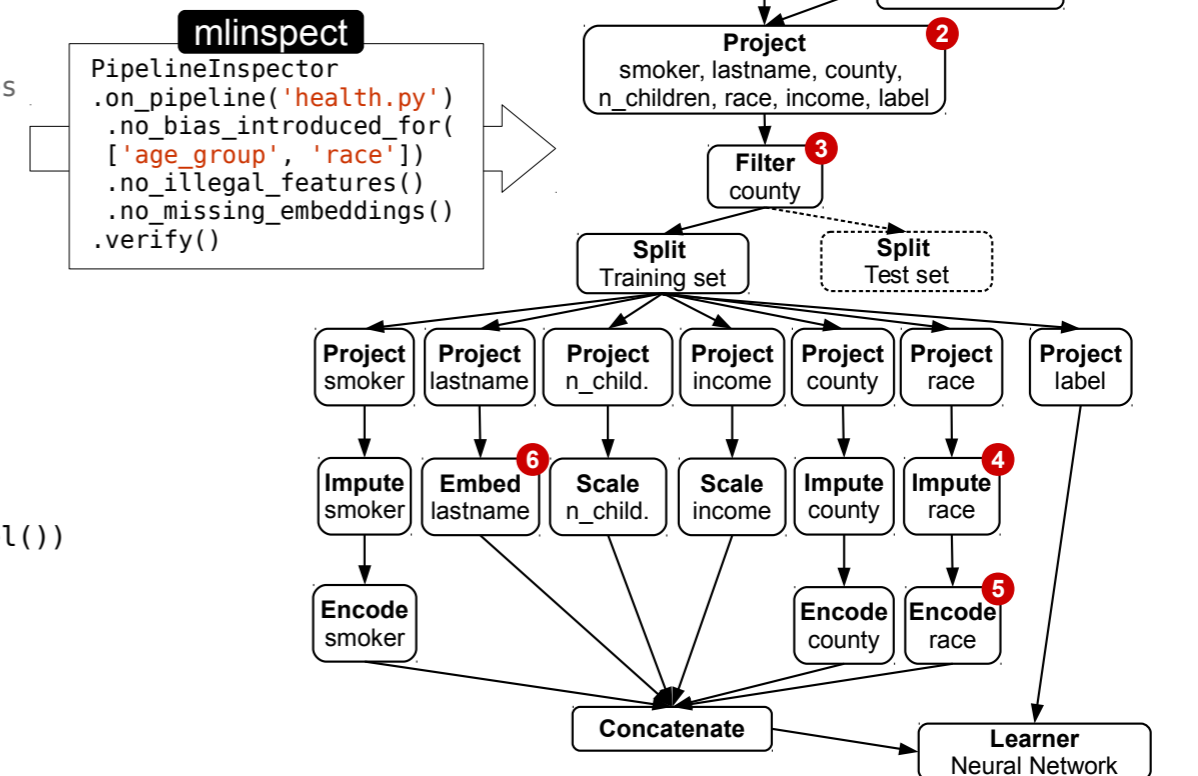
**Corresponding dataflow DAG for instrumentation, extracted by *mlinspect***

**Declarative inspection of preprocessing pipeline**

```
mlinspect
PipelineInspector
.on_pipeline('health.py')
.no_bias_introduced_for(
  ['age_group', 'race'])
.no_illegal_features()
.no_missing_embeddings()
.verify()
```

Data Source → Data Source → **Join** *on* ssn **1** → **Aggregate** *group by* age_group → **Join** *on* age_group → **Project** comp. / **Project** mean. → **Project** label → **Project** smoker, lastname, county, n_children, race, income, label **2** → **Filter** county **3** → **Split** Training set / **Split** Test set

**Project** smoker → **Impute** smoker → **Encode** smoker
**Project** lastname → **Embed** lastname **6**
**Project** n_child. → **Scale** n_child.
**Project** income → **Scale** income
**Project** county → **Impute** county → **Encode** county
**Project** race → **Impute** race **4** → **Encode** race **5**
**Project** label

→ **Concatenate** → **Learner** Neural Network

[Grafberger, Stoyanovich, Schelter (2021)]

center for responsible ai

# Automated Data Cleaning Can Hurt Fairness in ML-based Decision Making

Shubha Guha
s.guha@uva.nl
University of Amsterdam

Falaah Arif Khan
fa2161@nyu.edu
New York University

Julia Stoyanovich
stoyanovich@nyu.edu
New York University

Sebastian Schelter
s.schelter@uva.nl
University of Amsterdam

**ongoing work**

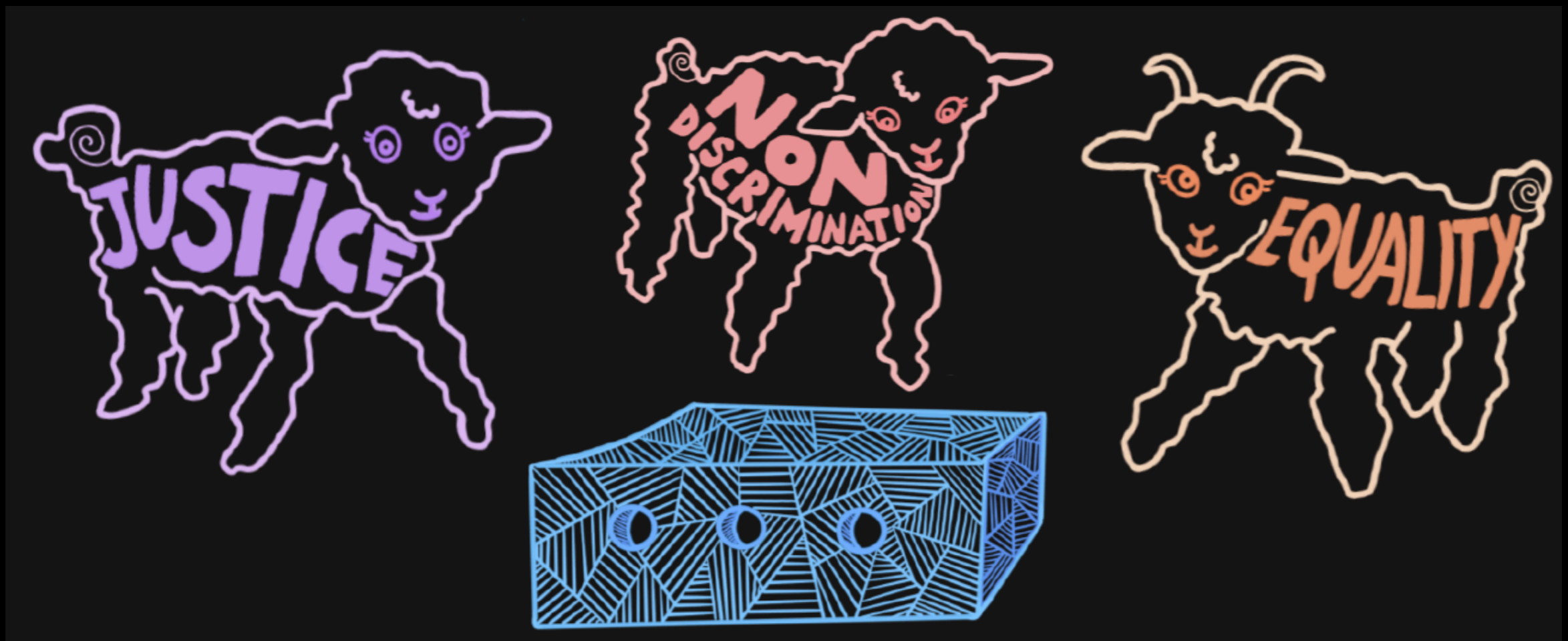| | auto-cleaning makes | | |
| model | fairness worse | fairness better | fairness & accuracy better |
|---|---|---|---|
| xgboost | 21.2% (45) | 10.8% (23) | 6.6% (14) |
| knn | 24.5% (52) | 13.7% (29) | 11.8% (25) |
| log-reg | 19.8% (42) | 12.3% (26) | 7.5% (16) |

TABLE V

IMPACT OF AUTO-CLEANING ON ACCURACY AND FAIRNESS FOR DIFFERENT ML MODELS ON 212 CONFIGURATIONS IN TOTAL. WE LIST CASES WHERE FAIRNESS GETS WORSE, FAIRNESS GETS BETTER, AND WHERE BOTH FAIRNESS AND ACCURACY GET BETTER. AUTO-CLEANING IS MORE LIKELY TO WORSEN THAN TO IMPROVE FAIRNESS ACROSS ALL MODELS.

**Emergent bias** arises in the context of use of a technical system

# THE WALL STREET JOURNAL.

## Hiring and AI: Let Job Candidates Know Why They Were Rejected

Labels that explain a hiring process that uses AI could allow job seekers to opt out if they object to the employer's data practices.
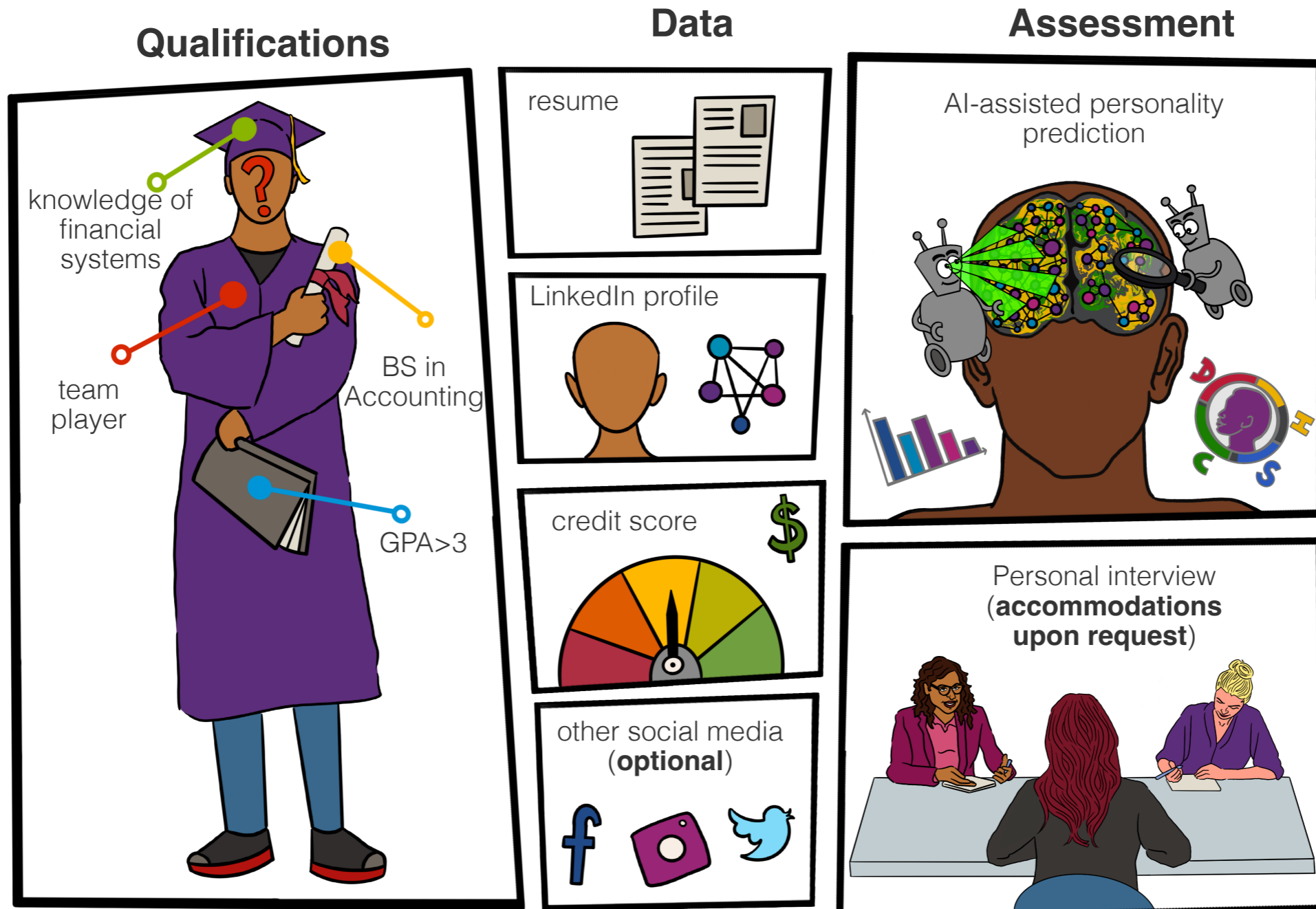PHOTO: ISTOCKPHOTO/GETTY IMAGES

*By Julia Stoyanovich*
Updated Sept. 22, 2021 11:00 am ET

Artificial-intelligence tools are seeing ever broader use in hiring. But this practice is also hotly criticized because we rarely understand how these tools select candidates, and whether the candidates they select are, in fact, better qualified than those who are rejected.

To help answer these crucial questions, **we should give job seekers more information about the hiring process and the decisions**. The solution I propose is a twist on something we see every day: **nutritional labels**. Specifically, job candidates would see simple, standardized labels that show the factors that go into the AI's decision.

r/ai center for responsible ai

https://dataresponsibly.github.io/we-are-ai/

# We are AI comics

r/ai center for responsible ai

# We are AI comics: in Spanish



Somos IA no. 1:
"¿QUÉ ES LA IA?"
© Julia Stoyanovich & Falaah Arif Khan (2022)

Somos IA no. 2:
APRENDER DE LOS DATOS
© Julia Stoyanovich & Falaah Arif Khan (2022)

Somos IA no. 3:
¿QUIÉN VIVE, QUIÉN MUERE, QUIÉN DECIDE?
© Julia Stoyanovich, Mona Sloane & Falaah Arif Khan (2022)

Somos IA no. 4:
TODO SOBRE ESE SESGO
© Julia Stoyanovich & Falaah Arif Khan (2022)

Somos IA no. 5:
SOMOS IA
© Julia Stoyanovich & Falaah Arif Khan (2022)

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

r/ai center for responsible ai

# Thank you!

**Julia Stoyanovich**
New York University
USA

**Serge Abiteboul**
Inria & ENS Paris
France

**Bill Howe**
University of Washington
USA

**H.V. Jagadish**
University of Michigan
USA

**Sebastian Schelter**
University of Amsterdam
The Netherlands