# Hearing #7 on Competition and Consumer Protection in the 21st Century

**Howard University**

**School of Law**

November 14, 2018

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

1

# Welcome

# We Will Be Starting Shortly

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

2

# Welcome and Introductory Remarks

**Bruce Hoffman**

Federal Trade Commission

Bureau of Competition

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

3

# Algorithmic Collusion

*Session moderated by:*

**Ellen Connelly**
Federal Trade Commission
Office of Policy Planning

**James Rhilinger**
Federal Trade Commission
Bureau of Competition

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

4

# Algorithmic Collusion

**Maurice E. Stucke**

University of Tennessee College of Law

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

5

# Algorithmic Collusion

**Ai Deng**

Bates White

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

6

# Algorithmic Collusion

**Kai-Uwe Kühn**

University of East Anglia

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

7

# Algorithmic Collusion

**Rosa M. Abrantes-Metz**

Global Economics Group

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

8

# Algorithmic Collusion

**Sonia Kuester Pfaffenroth**

Arnold & Porter

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

9

# Algorithmic Collusion

**Joseph E. Harrington, Jr.**

University of Pennsylvania

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

10

# Algorithmic Collusion

**Panel Discussion:**

Maurice E. Stucke, Ai Deng, Kai-Uwe Kühn,
Rosa M. Abrantes-Metz,
Sonia Kuester Pfaffenroth,
Joseph E. Harrington, Jr.,

**Moderators:** Ellen Connelly & James Rhilinger

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

11

# Break
# 10:45-11:00 am

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

12

# Framing Presentation (prerecorded)

## Michael I. Jordan

University of California, Berkeley

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

13

# Emerging Competition, Innovation, and Market Structure Questions Around Algorithms, Artificial Intelligence, and Predictive Analytics

*Session moderated by:*

**Brian O'Dea**
Federal Trade Commission
Bureau of Competition

**Nathan Wilson**
Federal Trade Commission
Bureau of Economics

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

14

# Emerging Competition, Innovation, and Market Structure Questions Around Algorithms, Artificial Intelligence, and Predictive Analytics

## Panel Discussion:

Robin Feldman, Joshua Gans,
Preston McAfee, Nicolas Petit

**Moderators:** Brian O'Dea & Nathan Wilson

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

15

# Facial Analysis Technology Warning Signs

## Joy Buolamwini

Algorithmic Justice League | MIT Media Lab

PhD, MIT Pending

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

16

# Automated Facial Analysis Tasks

# The Coded Gaze

Algorithmic bias creating exclusionary experiences discriminatory practices

# Silent Sweep: Over 117 Million US Adults in Face Surveillance Databases

One in two American adults is in a law enforcement face recognition network used in unregulated searches employing algorithms with unaudited accuracy.

The Perpetual Line Up
(Garvie , Bedoya, Frankle  2016)

# Real-World Impact

"In two cases [Scotland Yard Report], innocent women were matched with men."

- Ian Drury, The Daily Mail – May 15 2018

## 91% of South Wales Police's automated facial recognition matches

**wrongly identified innocent people**

**2,451 innocent people's** biometric photos taken and stored **without their knowledge**

# Expanding Use of Technology

# Potential Harms Index

| INDIVIDUAL HARMS | | COLLECTIVE SOCIAL HARMS |
|---|---|---|
| ILLEGAL DISCRIMINATION | UNFAIR PRACTICES | |
| HIRING | | LOSS OF OPPORTUNITY |
| EMPLOYMENT | | |
| INSURANCE & SOCIAL BENEFITS | | |
| HOUSING | | |
| EDUCATION | | |
| CREDIT | | ECONOMIC LOSS |
| DIFFERENTIAL PRICES OF GOODS | | |
| LOSS OF LIBERTY | | SOCIAL STIGMATIZATION |
| INCREASED SURVEILLANCE | | |
| STEREOTYPE REINFORCEMENT | | |
| DIGNATORY HARMS | | |

# Gender Shades

**Intersectional Accuracy Disparities in Commercial Gender Classification**

**230+ articles in 37+ countries on MIT Thesis Research findings**

**Buolamwini, J., Gebru, T. "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification." Proceedings of Machine Learning Research 81:1–15, 2018 Conference on Fairness, Accountability, and Transparency**

# Gold Standard Measures of Success Mislead



## Data is Destiny
**Does your data reflect the world?**

BENCHMARK SKEWS
## 80% PALE  75% MALE

# False Sense of Progress



**2014**
DEEPFACE

**97.35%**
ACCURACY ON
GOLD STANDARD
LFW BENCHMARK

(Taigman et al., 2014)

**GOLD STANDARD SKEWS**
Labeled Faces in The Wild
*Released in 2007*

*~77.5% Male*
*~83.5% White*

(Han and Jain, 2014)

# National Benchmarks Not Immune

## NIST 2015 IJB-A BENCHMARK

INTERSECTIONAL BREAKDOWN

**4.4% Darker Female**
20.2% Lighter  Female


**59.4% Lighter Male**
16% Darker Male

SINGLE AXIS

24.6% Female
**75.4% Male**

# Towards Better Evaluation



*PILOR PARLIAMENTS BENCHMARK*

*FIRST GENDER AND SKIN TYPE LABELED GENDER CLASSIFICATION BENCHMARK*

**54.4%** *Male*
**53.6%** *Lighter*

# Testing Commercial AI Systems

How accurate are systems from IBM, Microsoft, and Face++ at determining the gender of faces in inclusive benchmark?

# Overall Accuracy

Aggregate performance metrics can mask racial and gender bias

**93.7%**  **90%**  **87.9%**

**www.gendershades.org**

May 2017 PPB Results

# Gender Bias

## All companies perform better on men than women



| | FEMALE FACES | MALE FACES |
|---|---|---|
| Microsoft | 89.3% | 97.4% |
| FACE++ | 78.7% | 99.3% |
| IBM | 79.7% | 94.4% |

8-21% ERROR GAP

FEMALE          MALE

## www.gendershades.org

May 2017 PPB Results

# Skin Type ~ Racial Bias

## All companies perform better on whites than people of color



|  | DARKER FACES | LIGHTER FACES |
|---|---|---|
| Microsoft | 87.1% | 99.3% |
| FACE++ | 83.5% | 95.3% |
| IBM | 77.6% | 96.8% |

12-19% ERROR GAP

DARKER    LIGHTER

## www.gendershades.org

May 2017 PPB Results

# Intersectional Performance

| 94% | 79.2% | 100% | 98.3% |
|---|---|---|---|



| DARKER MALES | DARKER FEMALES | LIGHTER MALES | LIGHTER FEMALES |
|---|---|---|---|

May 2017 PPB Results

# Intersectional Performance

**99.3%**   **65.5%**   **99.2%**   **94.0%**



FACE++

*DARKER MALES*   *DARKER FEMALES*   *LIGHTER MALES*   *LIGHTER FEMALES*

May 2017 PPB Results

# Intersectional Performance

| 88% | 65.3% | 99.7% | 92.9% |
|---|---|---|---|



**DARKER MALES** | **DARKER FEMALES** | **LIGHTER MALES** | **LIGHTER FEMALES**

May 2017 PPB Results

# Further Disaggregation Uncovers Even Higher Error Rates



|  | TYPE I | TYPE II | TYPE III | TYPE IV | TYPE V | TYPE VI |
|---|---|---|---|---|---|---|
| Microsoft | 1.7% | 1.1% | 3.3% | 0% | 23.2% | 25.0% |
| FACE++ | 11.9% | 9.7% | 8.2% | 13.9% | 32.4% | **46.5%** |
| IBM | 5.1% | 7.4% | 8.2% | 8.3% | 33.3% | **46.8%** |

**Commercial Error Rates Per Skin Type on Female Labeled Faces in PPB

May 2017 PPB Results

# Company Responses to Gender and Racial Bias in Commercial AI Systems

IBM and Microsoft engaged researchers

All companies released new products within 7 months of receiving audit results

# Self-Reported Improvement

February 2018 Internal IBM Results

**98%**     **96.5%**     **99.8%**     **100%**



**DARKER MALES**     **DARKER FEMALES**     **LIGHTER MALES**     **LIGHTER FEMALES**

Self-Reported Results With .99 Treshhold

# External Follow-Up Evaluation

August 2018 PPB Results



**99.4%**    **83.0%**    **99.7%**    **97.6%**

**DARKER MALES**    **DARKER FEMALES**    **LIGHTER MALES**    **LIGHTER FEMALES**

Accuracy Determined Using Gender Label Returned By API

# Accuracy Doesn't Mitigate Abuse



Illustration: Sally Thurer for The Intercept/Getty Images

**IBM USED NYPD SURVEILLANCE FOOTAGE TO DEVELOP TECHNOLOGY THAT LETS POLICE SEARCH BY SKIN COLOR**

# Regulators Mitigate Abuse

**Gender Shades Research Supported Recommendations**
- Require Vendors of Facial Analysis Technology To:
  - Implement internal bias evaluation, mitigation, and reporting procedures
  - Regularly report performance on national benchmarks
  - Support independent evaluation from research community

- Require National Institute of Standards & Technology To:
  - Make public demographic and phenotypic composition of benchmarks
  - Report accessible intersectional performance metrics

# Regulators Mitigate Abuse

## Broader Considerations

- **Consent and Control:** Ensure consumers have meaningful opportunity to consent or refuse capture of face and ability to control use of face data – (Require companies like Facebook Provide Face Purge Option)

- **Transparency:** Require disclosure when facial analysis technology is in use and information about storage and use of face data

- **Due Process:** Provide mechanisms for redress and contestation of decisions made with or informed by facial analysis technology

- **Heightened Privacy:** Recognize that face images are identifying information, and enable processors to determine consumers' precise geolocation information

# For More Information Contact

## comms@ajlunited.org



Oprah Winfrey — amazon — appears to be male 76.5 %

Serena Williams — IBM WATSON — MALE 0.89

Michelle Obama — Microsoft — "a young man wearing a black shirt", "confidence": 0.7999446; "hairpiece", "confidence": 0.9350064

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

43

# Lunch
# 1:00-2:15 pm

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

44

# Fairness and Intelligibility in Machine Learning Systems

**Jenn Wortman Vaughan**

Microsoft Research

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

45

# The Age of AI



## NIPS Registrations

# New Challenges

**Online Ads for High-Paying Jobs Are Targeting Men More Than Women** New study uncovers gender bias

Do Google's 'unprofessional hair' results show it is racist?
Leigh Alexander

## When Algorithms Discriminate

When it Comes to Policing, Data Is Not Benign

The online world is shaped by forces beyond our control, determining the stories we read on Facebook, the people we meet on OkCupid and the search results we see on Google. Big data is used to make decisions about health care, employment, housing, education and policing.

**Amazon just showed us that 'unbiased' algorithms can be inadvertently racist**

Technology

**Google apologises for Photos app racist blunder**

© 1 July 2015 | Technology

# Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

O N A SPRING AFTERNOON IN 2014, Brisha Borden was running late to pick up her god-sister from school when she spotted an unlocked kid's blue Huffy bicycle and a silver Razor scooter. Borden

**Amazon Prime and the racist algorithms**

# Microsoft's AI Principles

**Fairness**

**Reliability & Safety**

**Privacy & Security**

**Inclusive-ness**

**Transparency**

**Accountability**

# FATE: Fairness, Accountability, Transparency, and Ethics in AI

Sensitive Uses of AI

AI Reliability and Safety

Human-AI Collab and Interaction

Fairness and Bias

Intelligibility & Explainability

Engineering Practices for AI

Human Attention & Cognition

# AETHER Committee

AI Ethics and Effects in Engineering and Research

# Partnership on AI
## to benefit people and society

aws · Google · DeepMind · Microsoft · facebook · Apple · IBM

FOUNDING PARTNERS

# What are machine learning and AI?

# AI

Computers doing
things that we
would normally
think of as
*intelligent*

# AI

Computers doing things that we would normally think of as *intelligent*

# MACHINE LEARNING

Systems that learn from DATA and EXPERIENCE instead of being explicitly programmed

# AI

Computers doing things that we would normally think of as *intelligent*

## MACHINE LEARNING

Systems that learn from DATA and EXPERIENCE instead of being explicitly programmed

NEURAL NETWORKS

# Types of Machine Learning

- **Supervised learning:** Use labeled data to learn a general rule mapping inputs to outputs

- **Unsupervised learning:** Identify hidden structure and patterns in data; cluster data points

- **Reinforcement learning:** Perform a task, such as driving a vehicle or playing a game, in a dynamic environment, learning through trial and error

# Why might a machine learning system be unfair?

# The Machine Learning Pipeline



Feedback → Task Definition → Dataset Construction → Model Definition → Training Process → Testing Process → Deployment → Feedback

# Task Definition

# Task Definition



(a) Three samples in criminal ID photo set $S_c$.

(b) Three samples in non-criminal ID photo set $S_n$

Figure 1. Sample ID photos in our data set.

(Wu and Zhang, 2016)

# Dataset Construction

# Data: Societal Bias

## 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.

# Data: Societal Bias



(Caliskan et al., 2017)

# Data: Societal Bias

# Data: Skewed Sample

| Gender Classifier | Darker Male | Darker Female | Lighter Male | Lighter Female | Largest Gap |
|---|---|---|---|---|---|
| Microsoft | 94.0% | 79.2% | 100% | 98.3% | 20.8% |
| FACE++ | 99.3% | 65.5% | 99.2% | 94.0% | 33.8% |
| IBM | 88.0% | 65.3% | 99.7% | 92.9% | 34.4% |

(Buolamwini and Gebru, 2018)

# Data: Labeler Bias



More States Opting To 'Robo-Grade' Student Essays By Computer

June 30, 2018 · 8:13 AM ET
Heard on **Weekend Edition Saturday**

TOVIA SMITH

# Model Definition

# Models are Mathematical Abstractions

price of house  =  $w_1$ * number of bedrooms

$+ w_2$ * number of bathrooms

$+ w_3$ * square feet

+ a little bit of noise

# Model: Assumptions



Artificial Intelligence Is Now Used to Predict Crime. But Is It Biased?

The software is supposed to make policing more fair and accountable. But critics say it still has a way to go.

# Training Process



Feedback → Task Definition → Dataset Construction → Model Definition → **Training Process** → Testing Process → Deployment → Feedback

# Training Process

price of house = $w_1$ * number of bedrooms

+ $w_2$ * number of bathrooms

+ $w_3$ * square feet

+ a little bit of noise

# Training Process

```python
if schedule['Lambda_SKCC'] <= self.total_iter:
    start = time.time()

    shp_SKCC[:] = np.outer(W_d_C, W_d_C)
    shp_SKCC[:, :, bool_diag_CC] = W_a_C * W_d_C
    shp_SKCC *= W_K[None, :, None, None]
    shp_SKCC *= W_S[:, None, None, None]
    post_shp_SKCC = shp_SKCC + Y_SKCC

    if mask.ndim == 2:
        mask_N...
        zeta_S... np.dot(Theta_NC.T, np.dot(mask_NN, Theta_NC))
        zeta... _TS.sum(axis=0)[:, None, None]
    else:
        mask_TNN = mask
        zeta_TNC = np.einsum('tij,jd->tid', mask_TNN, Theta_NC)
        zeta_TCC = np.einsum('tid,ic->tcd', zeta_TNC, Theta_NC)
        zeta_SCC = np.einsum('tcd,ts->scd', zeta_TCC, Psi_TS)
    post_rte_SKCC = d + zeta_SCC[:, None, :, :]
```

# Testing Process

# Testing: Metrics

**Translation tutorial:**
**21 fairness definitions and their politics**

**Arvind Narayanan**
**(Computer scientist, Princeton University)**

Computer scientists and statisticians have devised numerous mathematical criteria to define what it means for a classifier or a model to be fair. The proliferation of these definitions represents an attempt to make technical sense of the complex, shifting social understanding of fairness. Thus, these definitions are laden with values and politics, and seemingly technical discussions about mathematical definitions in fact implicate weighty normative questions. A core component of these technical discussions has been the discovery of trade-offs between different (mathematical) notions of fairness; these trade-offs deserve attention beyond the technical community.

# Testing: Metrics

|        | Unqualified | Qualified |
|--------|:-----------:|:---------:|
| Reject | TN          | FN        |
| Hire   | FP          | TP        |

# Testing: Metrics

|  | Unqualified | Qualified |
|---|---|---|
| Reject | TN | FN |
| Hire | FP | TP |

What is the probability that a woman is qualified given that you choose to hire her? What about a man?

Predictive parity requires (almost) equal values of

$$\frac{TP}{TP + FP}$$

# Testing: Metrics

| | Unqualified | Qualified |
|---|---|---|
| Reject | TN | FN |
| Hire | FP | TP |

What is the probability of hiring a woman if she is unqualified? What about a man?

False positive rate balance requires (almost) equal values of

$$\frac{FP}{FP + TN}$$

# Testing: Metrics

| | Unqualified | Qualified |
|---|---|---|
| Reject | TN | FN |
| Hire | FP | TP |

What is the probability of rejecting a woman if she is qualified? What about a man?

False negative rate balance requires (almost) equal values of

$$\frac{FN}{FN + TP}$$

# Testing: Metrics



Bernard Parker, left, was rated high risk; Dylan Fugett was rated low risk. (Josh Ritchie for ProPublica)

**Machine Bias**

There's software used across the country to predict future criminals.
And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

# Testing: Metrics

## RESPONSE TO PROPUBLICA: DEMONSTRATING ACCURACY EQUITY AND PREDICTIVE PARITY

The website ProPublica recently published a story that focused on the scientific validity of COMPAS, raising questions about racial bias. As a result of the article and the subsequent national attention that it garnered, Northpointe launched an in-depth analysis of the data samples used by ProPublica. Drawing from the results of our analysis of ProPublica's data, Northpointe unequivocally rejects the ProPublica conclusion of racial bias in the COMPAS risk scales.

Predictive modeling is a specialized field within statistics and the appropriate use and interpretation of valid predictive models require a solid understanding of the techniques and methodological nuances common to this type of work. Our detailed review of how ProPublica conducted their analysis revealed several statistical and technical errors such as misspecified regression models, mis-defined classification terms and measures of discrimination, the incorrect interpretation and use of model errors, and more. These errors led to a false conclusion of racial bias; we do not

# Testing: Metrics

**Monkey Cage**

## A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

By **Sam Corbett-Davies**, **Emma Pierson**, **Avi Feller** and **Sharad Goel**
October 17, 2016

(Kleinberg et al., 2016;
Chouldechova, 2017)

# Deployment



Feedback → Task Definition → Dataset Construction → Model Definition → Training Process → Testing Process → Deployment → Feedback

# Deployment: Context



**East Asian faces** / **Caucasian faces**

(Phillips et al., 2011)

# Feedback

# Feedback Loops

Use history of drug-crime reports and arrests to predict future crime locations…

⬇

More historic arrests in Black and Hispanic areas

⬇

More policing in these areas

⬇

More arrests in these areas

# So what can we do?

# Strategies to Mitigate Harms

- Prioritize fairness at every stage of the ML pipeline

- Think critically about implicit assumptions made at each stage

- Pay attention to potential biases in the data source and data preparation process

- Check if test data matches the deployment context

- Involve diverse stakeholders and gather multiple perspectives

- Acknowledge our mistakes and learn from them

# Transparency vs. Intelligibility

# What is Transparency?

- In policy circles, transparency represents two distinct ideas
  - People should be able to understand and monitor how AI systems work
  - Those who deploy AI systems should be honest and forthcoming about how and when they are being used

- In machine learning circles, the former is called "intelligibility" or "interpretability," and **literal transparency can work against it!**

# Transparency ≠ Intelligibility

- Exposing ML source code doesn't tell us much
- Exposing model internals can stop people from noticing when a model makes a mistake because of information overload



(Poursabzi-Sangdeh et al., 2018)

# Why intelligibility?

— Accountability: An applicant wants to know why she was denied a loan.

— Trust: A model deployed in a school predicts that a student is likely to drop out.  Knowing the factors relevant for the prediction could help his teacher decide whether to believe it and how to intervene.

— Bias assessment: A model matches candidates to jobs.  By understanding characteristics of the training data, an employer may see that female candidates are underrepresented, leading to potential bias.

— Robustness: A data scientist sees unexpected predictions from a model she has trained.  Knowing why these predictions were made could help her debug the model.

# Intelligibility via "Simple Models"



$$y = f_1(x_1) + \ldots + f_d(x_d)$$

## Point Systems
(Jung et al., 2017; Ustun & Rudin, 2015)

## Generalized Additive Models
(Lou, Caruana, et al., 2012&2013)

Classic methods: decision trees, rule lists (if-then-else), rule sets, sparse linear models, …

# Intelligibility via Post Hoc Explanations



## Simple Explanations of a Single Prediction
(e.g., Ribeiro et al., 2016; Lundberg and Lee, 2017)



## Simple Approximations of a Full Model
(e.g., Lakkaraju et al., 2017)

# Data Intelligibility: Datasheets for Datasets



(Gebru et al., 2018)

# Data Intelligibility: Datasheets for Datasets

- Questions cover dataset motivation, composition, collection process, pre-processing, distribution, maintenance, legal concerns, and ethical concerns

- Sample use cases:
  - Post with public datasets to inform potential users about the make-up and origin of the data
  - Include with a company's internal-use datasets to provide relevant information to future users from across the company

# No One-Size-Fits-All Solution

| | Audit a single prediction | Understand model globally | Make better decisions | Debug models | Assess bias | Inspire trust |
|---|---|---|---|---|---|---|
| CEOs | | | Approach A | | | |
| Data scientists | | | | Approach C | | |
| Lay people | | | | | | |
| Regulators | Approach B | | | | | |

# No One-Size-Fits-All Solution

– Why is the explanation needed? What is your goal?

– What is being explained? Prediction or whole system?

– To whom should the system be intelligible?

– Does the explainer have access to system internals?

– Does the explainer have access to the training data?

– What is the dimensionality or scale of the system?

– What type of data is used? Feature vectors? Text?

– Could giving away too much open up the system to manipulation?

– Could giving away too much reveal proprietary information?

# Takeaways

– There is no one-size-fits-all solution to fairness, transparency, or intelligibility

– These principles cannot be treated as afterthoughts; they must be considered at every stage of the machine learning pipeline

– Technology can be part of the solution, if used with care

– It is important to involve diverse stakeholders and gather multiple perspectives

– We should admit our mistakes and learn from them

# Thanks!

http://jennwv.com

jenn@microsoft.com

# Wrapping Up and Looking Ahead: Roundtable Discussion of Key Legal and Regulatory Questions in the Field

*Session moderated by:*

**Ellen Connelly**
Federal Trade Commission
Office of Policy Planning

**Benjamin Rossen**
Federal Trade Commission
Division of Privacy and Identity Protection

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

101

# Wrapping Up and Looking Ahead: Roundtable Discussion of Key Legal and Regulatory Questions in the Field

## Panel Discussion:

Justin Brookman, Pam Dixon,
Salil Mehra, Joshua New,
Nicol Turner-Lee

**Moderators:** Ellen Connelly & Benjamin Rossen

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

102

# Closing Remarks

## Danielle Holley-Walker

Howard University School of Law

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

103

# Thank You
# Join Us In December

Hearings on Competition and Consumer Protection in the 21st Century
An FTC-Howard University Law School Event | November 13-14, 2018 | ftc.gov/ftc-hearings | #ftchearings

104