

The background of the slide features a dramatic, large-scale industrial fire at night. The intense orange and yellow flames dominate the right side of the frame, casting a bright glow over dark structures that appear to be part of a factory or refinery. The fire is highly turbulent, with many small sparks and embers visible in the dark sky above.

# Machine Learning in Production Responsible ML Engineering



Nao Tokui

@naotokui\_en · [Follow](#)



"Success" and "Sadness", according to DALL-E 2.

(No cherry-picking)



4:00 AM · Aug 7, 2022



482

Reply

[Copy link](#)

[Read 21 replies](#)

# Changing directions...

## Fundamentals of Engineering AI-Enabled Systems

**Holistic system view:** AI and non-AI components, pipelines, stakeholders, environment interactions, feedback loops

### Requirements:

System and model goals  
User requirements  
Environment assumptions  
Quality beyond accuracy  
Measurement  
Risk analysis  
Planning for mistakes

### Architecture + design:

Modeling tradeoffs  
Deployment architecture  
Data science pipelines  
Telemetry, monitoring  
Anticipating evolution  
Big data processing  
Human-AI design

### Quality assurance:

Model testing  
Data quality  
QA automation  
Testing in production  
Infrastructure quality  
Debugging

### Operations:

Continuous deployment  
Contin. experimentation  
Configuration mgmt.  
Monitoring  
Versioning  
Big data  
DevOps, MLOps

**Teams and process:** Data science vs software eng. workflows, interdisciplinary teams, collaboration points, technical debt

## Responsible AI Engineering

Provenance,  
versioning,  
reproducibility

Safety

Security and  
privacy

Fairness

Interpretability  
and explainability

Transparency  
and trust

Ethics, governance, regulation, compliance, organizational culture

# Readings

R. Caplan, J. Donovan, L. Hanson, J. Matthews. "Algorithmic Accountability: A Primer", Data & Society (2018).

# Learning Goals

- Review the importance of ethical considerations in designing AI-enabled systems
- Recall basic strategies to reason about ethical challenges
- Diagnose potential ethical issues in a given system
- Understand the types of harm that can be caused by ML
- Understand the sources of bias in ML

# Why Fairness?

# Many interrelated issues:

\* Ethics \* Fairness \* Justice \* Discrimination \* Safety \* Privacy \* Security \* Transparency \* Accountability

*Each is a deep and nuanced research topic. We focus on survey of some key issues.*



*In 2015, Shkreli received widespread criticism [...] obtained the manufacturing license for the antiparasitic drug Daraprim and raised its price from USD 13.5 to 750 per pill [...] referred to by the media as "the most hated man in America" and "Pharma Bro". -- [Wikipedia](#)*

*"I could have raised it higher and made more profits for our shareholders. Which is my primary duty." -- Martin Shkreli*

## Speaker notes

Image source: [https://en.wikipedia.org/wiki/Martin\\_Shkreli#/media/File:Martin\\_Shkreli\\_2016.jpg](https://en.wikipedia.org/wiki/Martin_Shkreli#/media/File:Martin_Shkreli_2016.jpg)

# Terminology



**Legal** = in accordance to societal laws

- systematic body of rules governing society; set through government
- punishment for violation

**Ethical** = following moral principles of tradition, group, or individual

- branch of philosophy, science of a standard human conduct
- professional ethics = rules codified by professional organization
- no legal binding, no enforcement beyond "shame"
- high ethical standards may yield long term benefits through image and staff loyalty

# Big Disclaimer

Legality is obviously a locale-specific concern.

What is *ethical* (and how we know) is a very complicated question.

- Whether there exists ground-truth ethics is a point of philosophical debate.
- Often informed by context/culture/etc.

We adopt a generally US-centric perspective for much of this discussion.

- ...Because that's where we are.
- But given the global reach of software, tread with care.

Speaker notes

GDPR is an easy example



# With a few lines of code...

Developers have substantial power in shaping products, and software has substantial power over human lives.

Small design decisions can have substantial impact (safety, security, discrimination, ...) -- not always deliberate

Our view: We have both **legal & ethical** responsibilities to anticipate mistakes, think through their consequences, and build in mitigations!

# Example: Social Media



≡ *What is the (real) organizational objective of the company?*

# Optimizing for Organizational Objective

How do we maximize the user engagement? Examples:

- Infinite scroll: Encourage non-stop, continual use
- Personal recommendations: Suggest news feed to increase engagement
- Push notifications: Notify disengaged users to return to the app



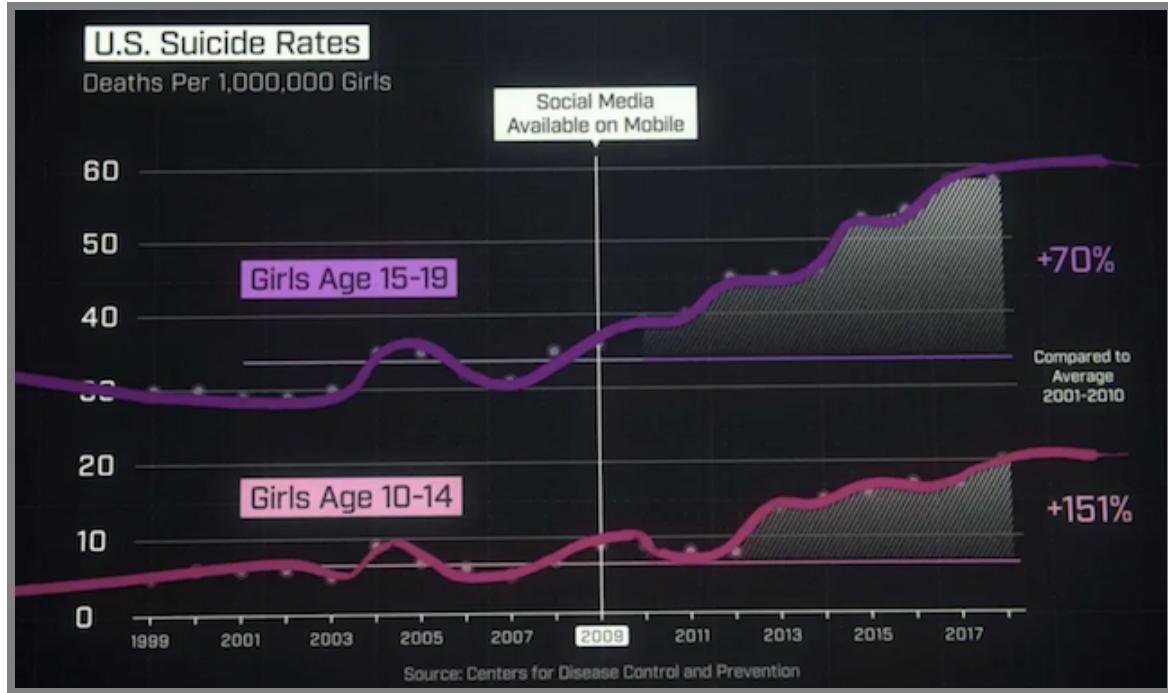
# Addiction



- 210M people worldwide addicted to social media
- 71% of Americans sleep next to a mobile device
- ~1000 people injured per day due to distracted driving (USA)

<https://www.flurry.com/blog/mobile-addicts-multiply-across-the-globe/>;  
[https://www.cdc.gov/motorvehiclesafety/Distracted\\_Driving/index.html](https://www.cdc.gov/motorvehiclesafety/Distracted_Driving/index.html)

# Mental Health



- 35% of US teenagers with low social-emotional well-being have been bullied on social media.
- 70% of teens feel excluded when using social media.



<https://leftronic.com/social-media-addiction-statistics>

# Disinformation & Polarization

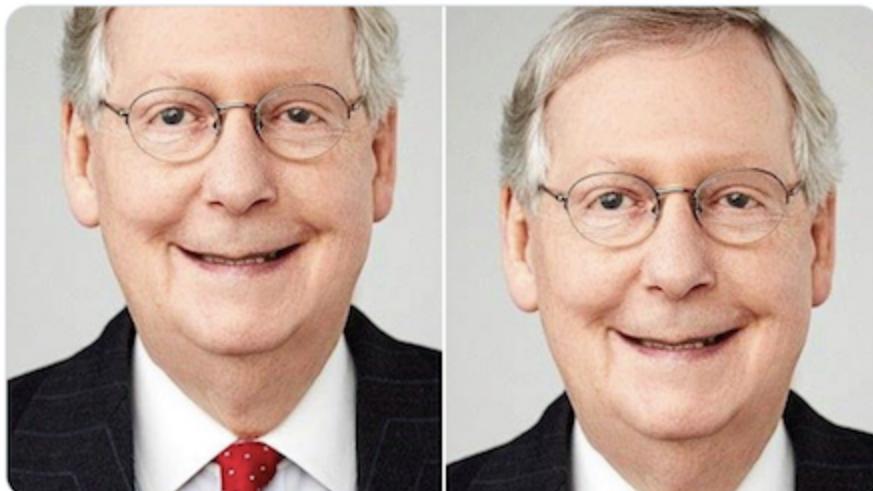


# Discrimination (as a side effect...)

 Tony "Abolish (Pol)ICE" Arcieri 🇺🇸  
@bascule

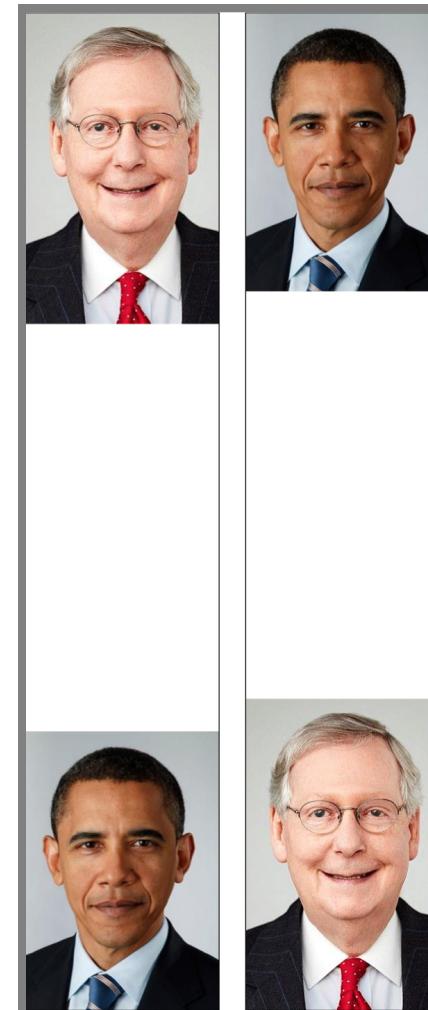
Trying a horrible experiment...

Which will the Twitter algorithm pick: Mitch McConnell or Barack Obama?



6:05 PM · Sep 19, 2020 · Twitter Web App

64K Retweets 16.5K Quote Tweets 198.3K Likes



# Who's to blame?

The screenshot shows a news article from The Independent. At the top left is the site's logo, 'INDEPENDENT'. To the right are links for 'Support us', 'Contribute' (in a red box), and 'Subscribe'. Below the header is a navigation bar with categories: NEWS, POLITICS, VOICES, SPORT, CULTURE, INDY/LIFE (in red), INDYBEST, VIDEO, DAILY EDITION, and CONVERSATIONS. The main title of the article is 'GOOGLE QUIETLY REMOVES 'DON'T BE EVIL' PREFACE FROM CODE OF CONDUCT', displayed in large, bold, black text on a white background. A subtitle below the main title reads 'Google employees resigned this month over the company's autonomous weapons project'. At the bottom of the article area, there is author information ('Anthony Cuthbertson | @ADCuthbertson | Monday 21 May 2018 12:21') and social sharing icons for bookmarking, Facebook, Twitter, and email.

Support us [Contribute](#) [Subscribe](#)

NEWS POLITICS VOICES SPORT CULTURE [INDY/LIFE](#) INDYBEST VIDEO DAILY EDITION CONVERSATIONS

## GOOGLE QUIETLY REMOVES 'DON'T BE EVIL' PREFACE FROM CODE OF CONDUCT

Google employees resigned this month over the company's autonomous weapons project

Anthony Cuthbertson | @ADCuthbertson | Monday 21 May 2018 12:21

*Are these companies intentionally trying to cause harm? If not, what are the root causes of the problem?*

# Liability?

*The software is provided “as is”, without warranty of any kind, express or implied, including but not limited to the warranties of merchantability, fitness for a particular purpose and noninfringement. In no event shall the authors or copyright holders be liable for any claim, damages or other liability, whether in an action of contract, tort or otherwise, arising from, out of or in connection with the software or the use or other dealings in the software.*

## Speaker notes

Software companies have usually gotten away with claiming no liability for their products



# Some Challenges

*Misalignment between organizational goals & societal values*

- Financial incentives often dominate other goals ("grow or die")

*Hardly any regulation*

- Often, little legal consequences for causing negative impact (with exceptions based on domain)
- Poor understanding of socio-technical systems by policy makers

*Engineering challenges, at system- & ML-level*

- Difficult to clearly define or measure ethical values
- Difficult to anticipate all possible usage contexts
- Difficult to anticipate impact of feedback loops
- Difficult to prevent malicious actors from abusing the system
- Difficult to interpret output of ML and make ethical decisions

**These problems have long existed, but are being rapidly exacerbated by the widespread use of ML**

There are ML-specific techniques/concerns with respect to these issues.

# Responsible Engineering Matters

Engineers have substantial power in shaping products and outcomes

Serious individual and societal harms possible from (a) negligence and  
(b) malicious designs

- Safety, mental health, weapons
- Security, privacy
- Manipulation, addiction, surveillance, polarization
- Job loss, deskilling
- Discrimination

# "I don't care about ethics, I just want to make money."

Regulations apply in many domains, including those where ML is "hot"

- Health care, finance, real estate

Bad PR can be bad for your bottom line.

DHH · Follow  
@dhh · Nov 7, 2019

The [@AppleCard](#) is such a fucking sexist program. My wife and I filed joint tax returns, live in a community-property state, and have been married for a long time. Yet Apple's black box algorithm thinks I deserve 20x the credit limit she does. No appeals work.

8:34 PM · Nov 7, 2019

23.6K · Reply · Copy link

# Responsible Engineering in this Course

Key areas of concern

- Fairness
- Safety
- Security and privacy
- Transparency and accountability

Technical infrastructure concepts

- Interpretability and explainability
- Versioning, provenance, reproducibility

# Fairness

# Dividing a Pie?

- Equal slices for everybody
- Bigger slices for active bakers
- Bigger slices for inexperienced/new members (e.g., children)
- Bigger slices for hungry people
- More pie for everybody, bake more



*(Not everybody contributed equally  
during baking, not everybody is  
≡ equally hungry)*

# What is fair?

*Fairness discourse asks questions about how to treat people and whether treating different groups of people differently is ethical. If two groups of people are systematically treated differently, this is often considered unfair.*

# Regulated domains (US)

- Credit (Equal Credit Opportunity Act)
- Education (Civil Rights Act of 1964; Education Amendments of 1972)
- Employment (Civil Rights Act of 1964)
- Housing (Fair Housing Act)
- 'Public Accommodation' (Civil Rights Act of 1964)

Extends to marketing and advertising; not limited to final decision

# Legally protected classes (US)

- Race ([Civil Rights Act of 1964](#))
- Religion ([Civil Rights Act of 1964](#))
- National origin ([Civil Rights Act of 1964](#))
- Sex, sexual orientation, and gender identity ([Equal Pay Act of 1963](#), [Civil Rights Act of 1964](#), and [Bostock v. Clayton](#))
- Age (40 and over, [Age Discrimination in Employment Act of 1967](#))
- Pregnancy ([Pregnancy Discrimination Act of 1978](#))
- Familial status (preference for or against having children, [Civil Rights Act of 1968](#))
- Disability status ([Rehabilitation Act of 1973](#); [Americans with Disabilities Act of 1990](#))
- Veteran status ([Vietnam Era Veterans' Readjustment Assistance Act of 1974](#); [Uniformed Services Employment and Reemployment Rights Act of 1994](#))
- Genetic information ([Genetic Information Nondiscrimination Act of 2008](#))

= [https://en.wikipedia.org/wiki/Protected\\_group](https://en.wikipedia.org/wiki/Protected_group)

# Common framing: Equality vs Equity vs Justice

# Equality vs Equity vs Justice



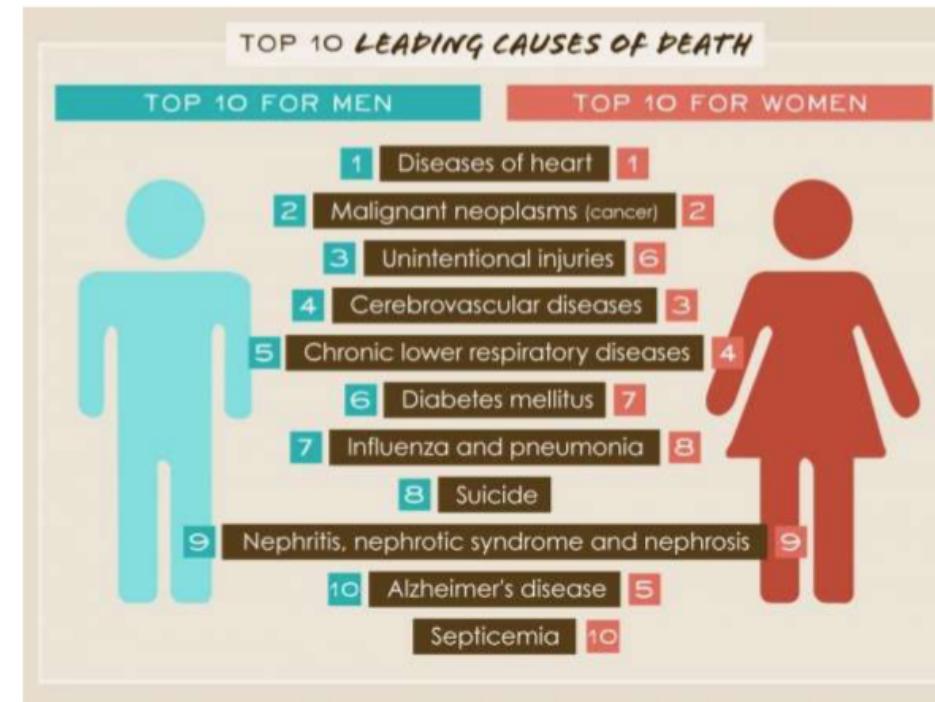
# Caveat: Something can be fair but still unethical (Thanos)!

# Not all discrimination is harmful



FEDERAL TRADE COMMISSION

Mortgage discrimination is against the law.



# Not all discrimination is harmful



- Discrimination is a **domain-specific** concept
  - ML models discriminate based on input data by construction.
  - There are real differences between two groups, it might not be fair to ignore them
  - The problem is *unjustified* differentiation; i.e., discriminating on factors that should not matter

# Fairness vs. bias vs. harm

Fairness is best understood as a **societal or cultural concept**.

Bias, in discussing ML, can be understood as a **technological or algorithmic concept**; it is often discussed in terms of its negative effects.

- Whether bias is harmful or unfair is not something that can be decided algorithmically.

Useful definition/framework defines algorithmic bias as "a skew that produces a type of harm."

# Types of Harm on Society

**Harms of allocation:** Withhold opportunities or resources

**Harms of representation:** Reinforce stereotypes, subordination along the lines of identity

# Harms of Allocation

- Withhold opportunities or resources
- Poor quality of service, degraded user experience for certain groups



# Harms of Representation

- Over/under-representation of certain groups in organizations
- Reinforcement of stereotypes (e.g. Black community & criminality)



"Racially identifying names" change the ads you get -- names commonly associated with Black individuals were more likely to trigger ads that suggested a criminal background check.

≡ *Discrimination in Online Ad Delivery, Latanya Sweeney, SSRN (2013).*

# Identifying (co-occurring) harms

	Allocation of resources	Quality of Service	Stereotyping	Denigration	Over- / Under-Representation
Hiring system does not rank women as highly as men for technical jobs	x	x	x		x
Photo management program labels image of black people as “gorillas”		x		x	
Image searches for “CEO” yield only photos of white men on first page			x		x

- Multiple types of harms can be caused by a product!
- Think about your system objectives & identify potential harms.

 Challenges of incorporating algorithmic fairness into practice, FAT\* Tutorial (2019). \*

# Role of Requirements Engineering

- Identify system goals
- Identify legal constraints
- Identify stakeholders and fairness concerns
- Analyze risks with regard to discrimination and fairness
- Analyze possible feedback loops (world vs machine)
- Negotiate tradeoffs with stakeholders
- Set requirements/constraints for data and model
- Plan mitigations in the system (beyond the model)
- Design incident response plan
- Set expectations for offline and online assurance and monitoring

# Sources of Bias

# Where does the bias come from?

The image shows two side-by-side screenshots of the Google Translate interface, demonstrating how machine learning models can exhibit gender bias.

**Top Screenshot (English to Turkish):**

- Source text: "He is a nurse  
She is a doctor"
- Target text: "O bir hemşire  
O bir doktor"
- Annotations: The word "he" is highlighted in blue, and the word "she" is highlighted in pink.

**Bottom Screenshot (Turkish to English):**

- Source text: "O bir hemşire  
O bir doktor"
- Target text: "She is a nurse  
He is a doctor" (Note: The target text is swapped compared to the source)
- Annotations: The word "she" is highlighted in blue, and the word "he" is highlighted in pink.

In both cases, the model appears to associate "he" with "nurse" and "she" with "doctor", which is a well-known gender bias in language corpora used for training.

= Semantics derived automatically from language corpora contain human-like biases, Caliskan et al., Science (2017).

# Where does the bias come from?

The image shows a screenshot of the Microsoft Translator interface, which displays four panels of text translation between English and Turkish.

- Top Left Panel (English to Turkish):** Shows the English sentence "He is a nurse.  
She is a doctor." Translated into Turkish as "O bir hemşire.  
O bir doktor." This panel has a character limit of 31/5000.
- Top Right Panel (Turkish to English):** Shows the Turkish sentence "O bir hemşire.  
O bir doktor." Translated back into English as "She's a nurse.  
He's a doctor." This panel also has a character limit of 31/5000.
- Bottom Left Panel (Turkish to English):** Shows the Turkish sentence "O bir hemşire.  
O bir doktor." Translated into English as "She's a nurse.  
He's a doctor." This panel has a character limit of 28/5000.
- Bottom Right Panel (English to Turkish):** Shows the English sentence "She's a nurse.  
He's a doctor." Translated into Turkish as "O bir hemşire.  
O bir doktor." This panel has a character limit of 28/5000.

The interface includes a Microsoft logo at the top left, a search bar at the top right, and navigation links for Translator, Text, Conversation, Apps, For business, and Help. A "Sign in" button is also visible at the top right.

# Sources of Bias

- Historial bias
- Tainted examples
- Limited features
- Skewed sample
- Sample size disparity
- Proxies



*Big Data's Disparate Impact*, Barocas & Selbst California Law Review (2016).

# Historical Bias

*Data reflects past biases, not intended outcomes*



*Should the algorithm reflect reality?*

## Speaker notes

"An example of this type of bias can be found in a 2018 image search result where searching for women CEOs ultimately resulted in fewer female CEO images due to the fact that only 5% of Fortune 500 CEOs were woman—which would cause the search results to be biased towards male CEOs. These search results were of course reflecting the reality, but whether or not the search algorithms should reflect this reality is an issue worth considering."



# What is reality?

## Speaker notes

at least as many women play games casually as men if not more. But someone searching for "gamer" probably has an image in their head. Should we be accurate with respect to reality? Or accurate with respect to what the person is searching for?



# Correcting Historical Bias

Fix the system, not just the model



A Scalable Approach to Reducing Gender Bias in Google Translate

# Correcting Historical Bias

*"Big Data processes codify the past. They do not invent the future. Doing that requires moral imagination, and that's something only humans can provide. " -- Cathy O'Neil in [Weapons of Math Destruction](#)*

*"Through user studies, the [image search] team learned that many users were uncomfortable with the idea of the company “manipulating” search results, viewing this behavior as unethical." -- observation from interviews by Ken Holstein*

# Correcting Historical Bias

*"Big Data processes codify the past. They do not invent the future. Doing that requires moral imagination, and that's something only humans can provide. " -- Cathy O'Neil in [Weapons of Math Destruction](#)*

*"Through user studies, the [image search] team learned that many users were uncomfortable with the idea of the company “manipulating” search results, viewing this behavior as unethical." -- observation from interviews by Ken Holstein*

# Correcting Historical Bias is Hard, even if you want to



≡ Google's 'CEO' image search gender bias hasn't really been fixed

# Tainted Labels

*Bias in dataset labels assigned (directly or indirectly) by humans*

The screenshot shows a news article from Vice Media. At the top left, there are three colored tabs: 'TECH' in pink, 'AMAZON' in blue, and 'ARTIFICIAL INTELLIGENCE' in red. The main title of the article is 'Amazon reportedly scraps internal AI recruiting tool that was biased against women'. Below the title is a subtitle: 'The secret program penalized applications that contained the word "women's"'. At the bottom left, it says 'By James Vincent | Oct 10, 2018, 7:09am EDT'.

**Amazon reportedly scraps internal AI recruiting tool that was biased against women**

*The secret program penalized applications that contained the word “women’s”*

By James Vincent | Oct 10, 2018, 7:09am EDT

Example: Hiring decision dataset -- labels assigned by (possibly biased) experts or derived from past (possibly biased) hiring decisions

# Limited Features

*Features that are less informative/reliable for certain subpopulations*



- Graduate admissions: Letters of recommendation equally reliable for international applicants?
- Employee performance review: "Leave of absence" acceptable feature if parental leave is gender skewed?

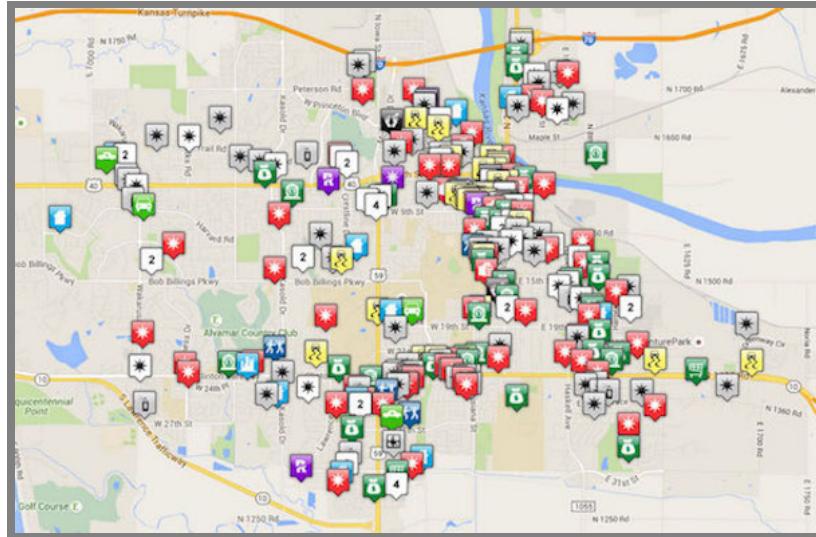
## Speaker notes

Decisions may be based on features that are predictive and accurate for a large part of the target distribution, but not so for some other parts of the distribution. For example, a system ranking applications for graduate school admissions may heavily rely on letters of recommendation and be well calibrated for applicants who can request letters from mentors familiar with the culture and jargon of such letters in the US, but may work poorly for international applicants from countries where such letters are not common or where such letters express support with different jargon. To reduce bias, we should be carefully reviewing all features and analyze whether they may be less predictive for certain subpopulations.



# Skewed Sample

*Bias in how and what data is collected*



Crime prediction: Where to analyze crime? What is considered crime?  
Actually a random/representative sample?

Raw data is an oxymoron

# Sample Size Disparity

*Limited training data for some subpopulations*



- Biased sampling process: "Shirley Card" used for Kodak color calibration, using mostly Caucasian models
- Small subpopulations: Sikhs small minority in US (0.2%) barely represented in a random sample

# Sample Size Disparity

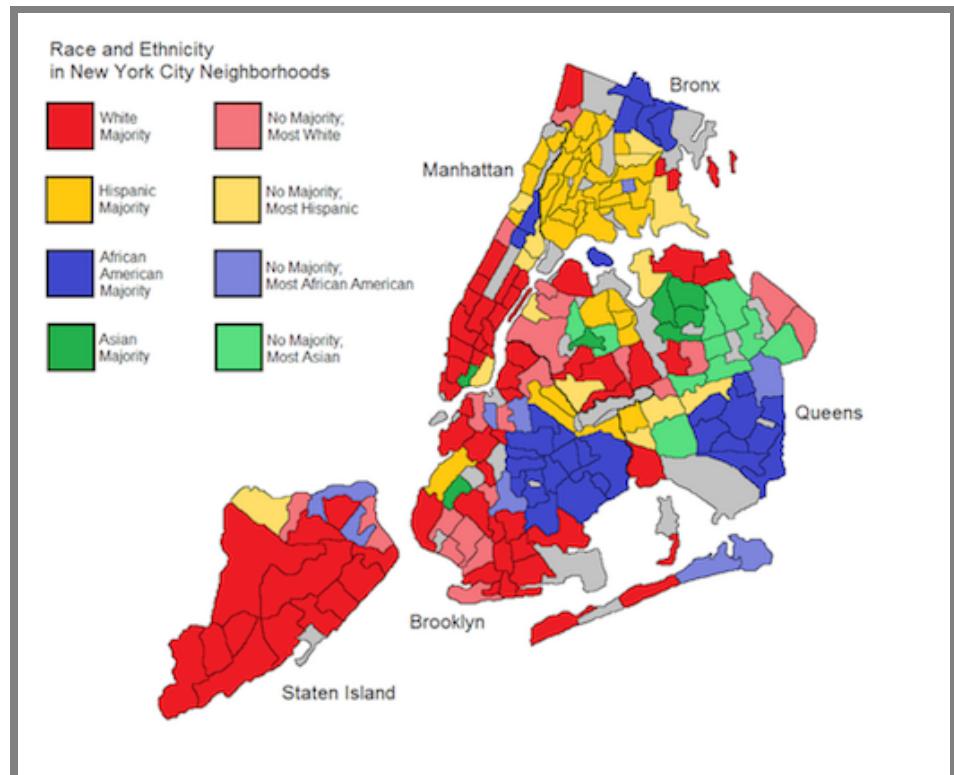
Without intervention:

- Models biased toward populations more represented in target distribution (e.g., Caucasian skin tones)
- ... biased towards population that are easier to sample (e.g., people self-selecting to post to Instagram)
- ... may ignore small minority populations as noise

Typically requires deliberate sampling strategy, intentional oversampling

# Proxies

*Features correlate with protected attribute, remain after removal*



- Example: Neighborhood as a proxy for race
- Extracurricular activities as proxy for gender and social class (e.g., “cheerleading”, “peer-mentor for ...”, “sailing team”, “classical music”)

# Feedback Loops reinforce Bias



*"Big Data processes codify the past. They do not invent the future. Doing that requires moral imagination, and that's something only humans can provide. " -- Cathy O'Neil in [Weapons of Math Destruction](#)*

# Breakout: College Admission



Scenario: Evaluate applications & identify students likely to succeed

Features: GPA, GRE/SAT, gender, race, undergrad institute, alumni connections, household income, hometown, transcript, etc.

# Breakout: College Admission

Scenario: Evaluate applications & identify students who are likely to succeed

Features: GPA, GRE/SAT, gender, race, undergrad institute, alumni connections, household income, hometown, transcript, etc.

As a group, post to #lecture tagging members:

- **Possible harms:** Allocation of resources? Quality of service?  
Stereotyping? Denigration? Over-/Under-representation?
- **Sources of bias:** Skewed sample? Tainted labels? Historical bias?  
Limited features? Sample size disparity? Proxies?

# Next lectures

1. Measuring and Improving Fairness at the Model Level
2. Fairness is a System-Wide Concern

# Summary

- Many interrelated issues: ethics, fairness, justice, safety, security, ...
- Both legal & ethical dimensions
- Challenges with developing ethical systems / developing systems responsibly
- Large potential for damage: Harm of allocation & harm of representation
- Sources of bias in ML: Skewed sample, tainted labels, limited features, sample size, disparity, proxies

# Further Readings

- O’Neil, Cathy. [Weapons of math destruction: How big data increases inequality and threatens democracy](#). Crown Publishing, 2017.
- Barocas, Solon, and Andrew D. Selbst. “[Big data’s disparate impact](#).” Calif. L. Rev. 104 (2016): 671.
- Mehrabi, Ninareh, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. “[A survey on bias and fairness in machine learning](#).” ACM Computing Surveys (CSUR) 54, no. 6 (2021): 1–35.
- Bietti, Elettra. “[From ethics washing to ethics bashing: a view on tech ethics from within moral philosophy](#).” In Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, pp. 210–219. 2020.

