

Week 7: Ethics

Charlottesville Handout

E t h i c s i S i t
E t h i c s i n a l g o r i t h m e a n ?

ETHICS in algorithmic hiring mean?

"Ethics in algorithmic hiring means "Dapt aai nettsmii ag spatiècsta uarbee tt hea we only need to consider the ethicality of the hiring algorithm."

Last week we discussed data privacy in the context of a lot of our conversations at this data call. • How data is consumed and services

- How data is reused

I sit ethically for data about us to be collected without our knowledge?

Companies are hiring diverse.

Data ethics... you mean machine learning?

Machine learning, big data and AI are all tremendous ways of dealing with sensitive and private data.

But data ethics is important every time we're dealing with sensitive and/or private data.

Even if we've got survey data about attitudes to green spaces during the pandemic.

In this course when we talk about data ethics we're also talking about the ethics of algorithms.

E th i c s & Mor a l

Ethics & Moral Philosophy

Moral philosophy and the history of thought is fascinating.

We're just going to hit some highlights.

R i g h t o p r i v a t e

R i g h t o p r i v a c y

But still be took at a longer distance a little bit more

S e l f - d r i v i n g

Sessions - driving change

The second continent of "smart" technology discussion focuses on what follows and finance the teach as follows. Who does the automation respond to? The responder under the heading car priority control of an automatic system that controls a steering wheel driving and breaking and driver does not have a different attitude towards the car can design the automated driving system and resume control.

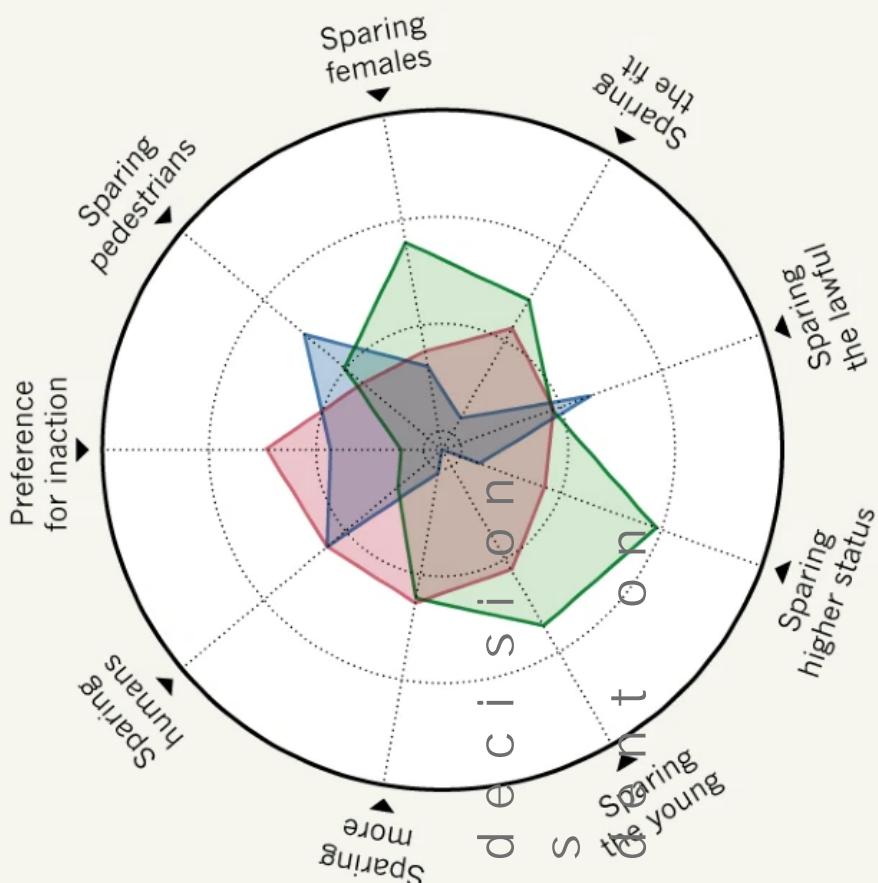
What techniques does the addition of self-reinforcement to net works raise?

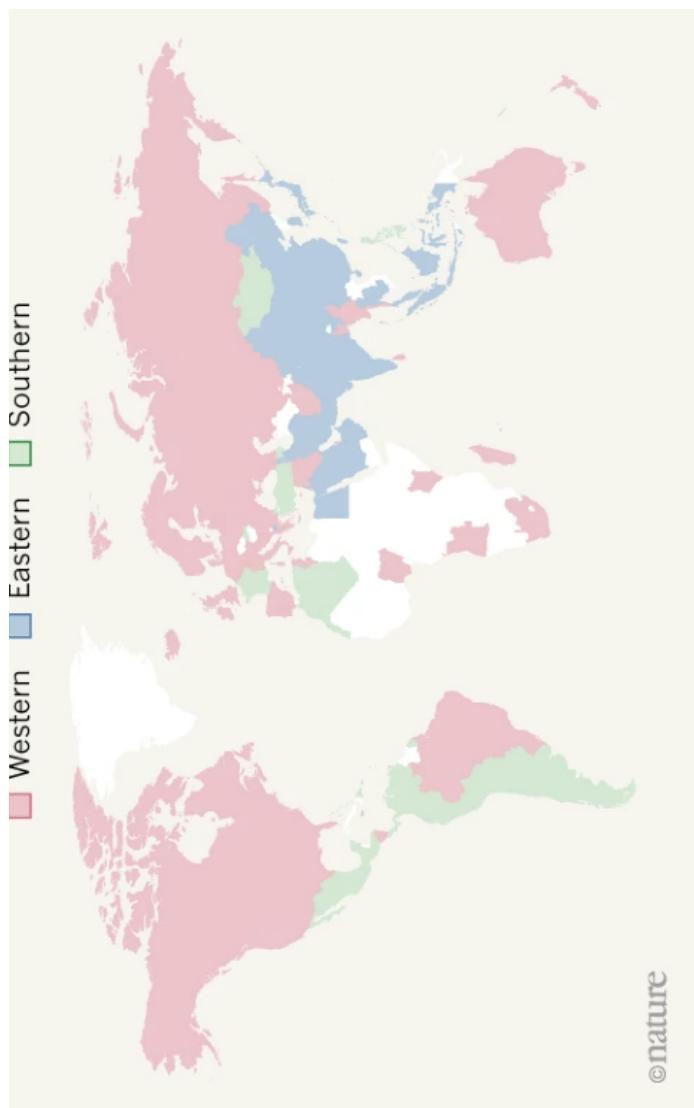
Ses - f - d r i v i n g C a r s (-)

Awarded to the university of people in 2018 to tomorrow's emma s autonomous vehicles.

Atitudenes vary by region of partisanship and emigration status, controlling for sex.

Should the
making of
vehicles
location?





This data visualization is derived from Awad et al.³ but appears in Maxmen 2018⁴.

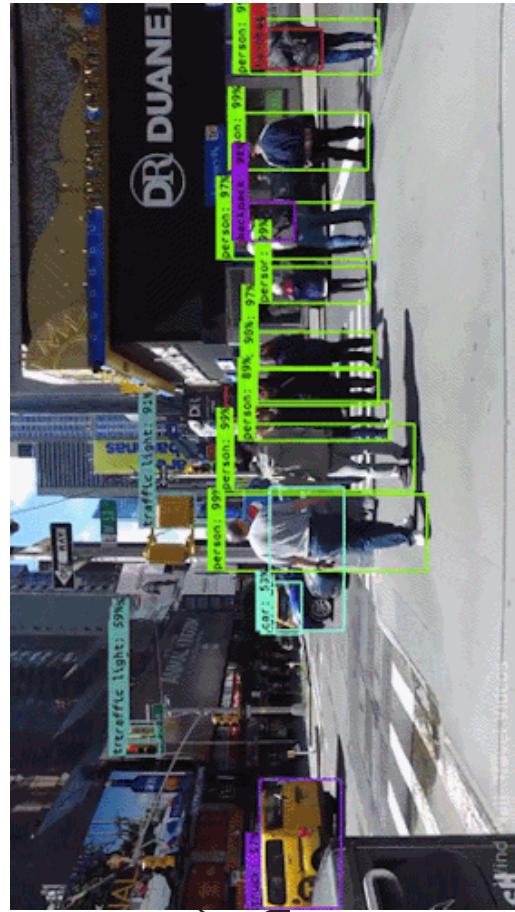
Ses - f - d r i v i n g C a r s (-)

Reiliabiliteenhicline
Tmai tsiosntsaestisaiuctoisnsc
easidunbiassendunbiandunbi
fintervperreytweheerreenin202

As more autonomous vehicles in real-world environments make mistakes, we throw away and reuse information from sensors to refine decisions on the fly. But this important work requires more information about autonomous vehicles than what's available in the literature. In fact, more information is needed to train drivers in real-world situations. This is where our research comes in.

Self-driving Cars (I)

Self-driving cars often use sensors for sensing and analysis at a greater what environment.



These considerable developments put about two common technologies: LiDAR and vision.

- LiDAR continuously measures a 3D visual map. This is used to detect objects around the car

S e l - f - d r i v i n g
S e l - f - d r i v i n g c a r s (—)

As we question the ethicality of autonomous vehicles making under-shot decisions particularly during a crash, we ask the following?

Does the vehicle correctly handle (cf) situations around a TnspA (ATe) and a DrAtnS (pgrgnb)?

Thankfully LiDAR and providers are providing answers to our questions from a pedestrian perspective.

We can see the autonomy mostly relies on inference to identify a pedestrian.

| | | | Classification and Path Prediction ^a | Vehicle and System Actions ^b |
|-------------------|--------------------------|-------------|--|---|
| | Time to Impact (seconds) | Speed (mph) | | |
| -9.9 | | 35.1 | -- | Vehicle begins to accelerate from 35 mph in response to increased speed limit. |
| -5.8 | | 44.1 | -- | Vehicle reaches 44 mph. |
| -5.6 | | 44.3 | Classification: Vehicle—by radar Path prediction: <i>None</i> ; not on path of SUV | Radar makes first detection of pedestrian (classified as vehicle) and estimates speed. |
| -5.2 | | 44.6 | Classification: Other—by lidar Path prediction: <i>Static</i> ; not on path of SUV | Lidar detects unknown object. Object is considered new, tracking history is unavailable, and velocity cannot be determined. ADS predicts object's path as static. |
| -4.2 | | 44.8 | Classification: Vehicle—by lidar Path prediction: <i>Static</i> ; not on path of SUV | Lidar classifies detected object as vehicle; this is a changed classification of object and without a tracking history. ADS predicts object's path as static. |
| -3.9 ^c | | 44.8 | Classification: Vehicle—by lidar Path prediction: Left through lane (next to SUV); not on path of SUV | Lidar retains classification <i>vehicle</i> . Based on tracking history and assigned goal, ADS predicts object's path as traveling in left through lane. |
| -3.8 to -2.7 | | 44.7 | Classification: alternates between vehicle and other—by lidar Path prediction: alternates between static and left through lane; neither considered on path of SUV | Object's classification alternates several times between vehicle and other. At each change, tracking history is unavailable; ADS predicts object's path as static. When detected object's classification remains same, ADS predicts path as traveling in left through lane. |
| -2.6 | | 44.6 | Classification: Bicycle—by lidar Path prediction: <i>Static</i> ; not on path of SUV | Lidar classifies detected object as <i>bicycle</i> ; this is a changed classification of object and object is without a tracking history. ADS predicts bicycle's path as static. |
| -2.5 | | 44.6 | Classification: Bicycle—by lidar Path prediction: Left through lane (next to SUV); not on path of SUV | Lidar retains <i>bicycle</i> classification; based on tracking history and assigned goal, ADS predicts bicycle's path as travelling in left through lane. |

Fairness, Accuracy Transparency

Fairness, Accuracy Training Parity

Accountability

Faînăsește

Fairness:

5 (Or 6)

Frustration, a tiny fingerally the source of this! But this is it since at least 2014

- Proxies for this! It might be a frequent user since at least 2014
- Skewed sample
- Trained on different sources of common sources of bias in training datasets.
- Broadly in the literature there are biases in training datasets.
- These biases result in our fair treatment of individuals / groups by law enforcement agencies, including police, or addenda to a contract (or addenda to a contract).
- Refusal of service
- More expensive services than less expensive services
- Reduced range of services

When these services include health care the findings can be life threatening.

Fairness:

Proxies are the easiest to identify and define

When training our algorithm we want to remember genealogical relationships - because we consider a language or flag a page a specific

However, many porters vary a lot in fairness. Can you think of any?

- Systemic racism means that income, neighborhood for race.

"Redlining is the practice of generally abusing it. So it's

- Facebook friends can be a strong proxy

Fairness: Limited

This is a harder bias than ~~directive money~~ (~~fairness~~ ~~money~~)
down examples

Limited features is a consequence of having
specific combinations of sensitive attributes.

Individuals in these groups will treat
than other groups.

This article gives a theoretical example:
unintended biased training - data - based 3347

Limited features is a pre-request for some

Fairness:

Skewed Sample Size Parities

I st' kind of unfair this is collected together
exampler, sample size disparities

Skewed samples is a bias in your algorithm
result in a heterogeneous model we do averaging maybe
feed back purposes do reas^{d2} paper Ensign 2013 also
really well.

Given historical crime incidents data for
patrols often areas to detect crime
....

Since such discovered incidents only occur
sent to by the predictive policing algorithm
bias to be compounded, causing away from

Fairness : Skewed Search

Another really great example with an easy issue with Professors like Avzee rp oewte¹³ ra | 2014

In 2009 Google published street address(GaFbTo) update¹⁴ that used search results that so. predict

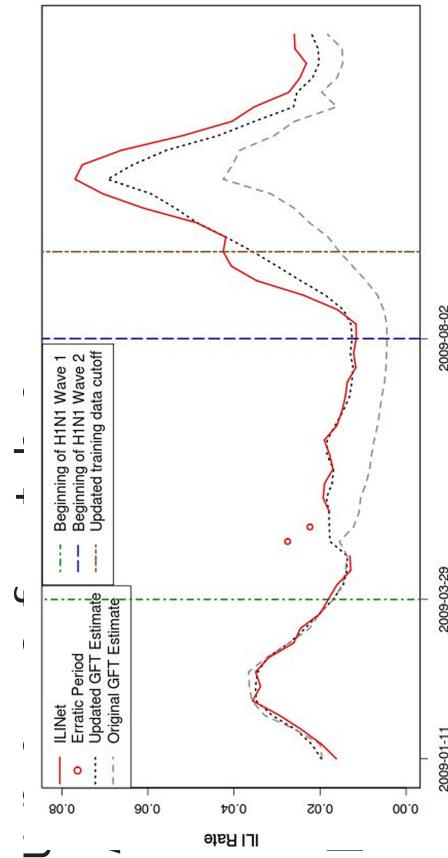
The theory being that ill people search for detected.

However whole thinking is some what questiona

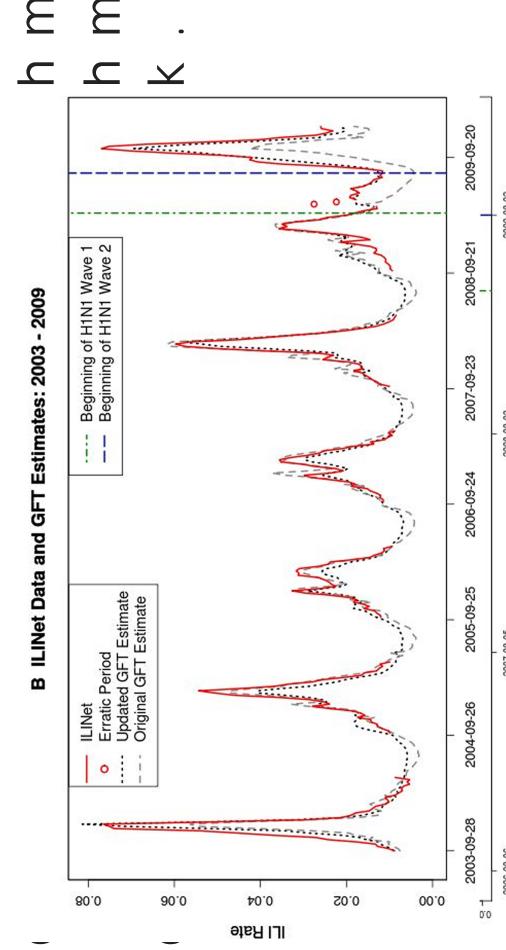
GFT has never documented the 45 searching been released a lot earlier mislead in Source LZCZerrer et¹³ a | 2014

Fairness : Skewed

The GFT didn't detect H1N1 during wave of H1N1 2015 and the service was the 2009.

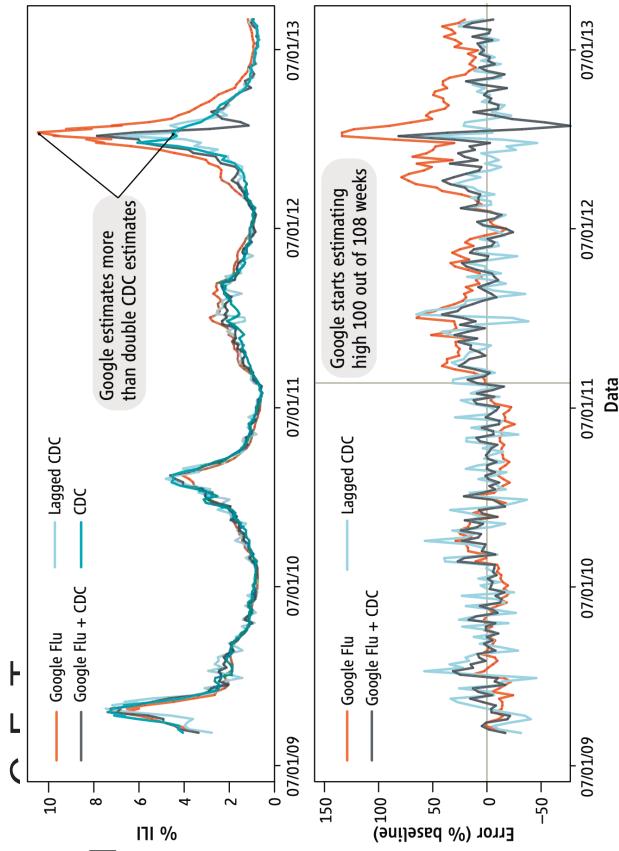


Google quickly updated I later in 2009 and the a better job of predicting a better period.



Fairness : Skewed Sampling

However from March to June
consistent flu volume
and became increasing
Lazer¹³ reported two reasons
• Big data hubris
Was there a global flu detector?



Source: Lazear et al

- Algorithm dynamically

Google has many competing things at the same time (or worse)

Fairness in AI

We need to talk about learning.

Do you know at the end supervised and unsupervised categories.

In supervised learning categories and assess these categories.

The algorithm will also belong to one

Thirst 'a l | well | and |
sometimes | i k | et | id | a
survival dataset

... but often our cause
absolute, e.g.

- Google uses AI has a
- Rating's Alirik & their ho
accident
- Rating a "good job

Lest 'l ook into the
example more

Fairness in Examples

Algoritmos are used to build high quality systems to address some of the challenges of the current literature. Fraiján et al. review literature reviews in 2021.

These algorithms are often used for matching resumes to job offers. They are also used for hiring decisions.

- Who gets hired: And who is rejected? For being biased and allowing more interesting people to do it.

This could not possibly create issues.

Fairness: Sample Size Discrepancy

Unfortunately hyperfected y balance [unintended] unbalanced categories can have bias. Same sample sizes.

Sample size disparity Tehwei smotss tw hwei nd el y c i t e d complete balance aetxaasneptl se hoafv et hi s i g - th disparity in the size **hotft possu: b/g/reonu. p*ssi*. K i p*e.d i a*.**

This disparity results in unfair behaviour between a population - real world datasets.

Fairness : Sample Size

In 2009, I implemented a mechanism to detect policy for Google+. I ast name patterns found their names rejected.

There aficionado vacay issues with this issue's vulnerability. This issue's was dealing same attack sets having same pattern here were fewer instances of higher names rejected.

But users also found "not real".

I recommended reading the **Fairness** paper. **Nameless** - What Happened to Programming Languages?

S o m e
m o r e
c a s e

...

"Man is to computer programmer as (I)

In 2016 Bo ²⁰ purchased it a table for example, if you know the exact word vocabulary of a algorithm for ~~is~~ ²¹ given words and their definitions which should look like word2vec a patent filed a artificial gorithm when trained on a given text find similar synonymy nouns words input.

Google provides a useful interface
however the system works.

"Man is to computer programmer as
()

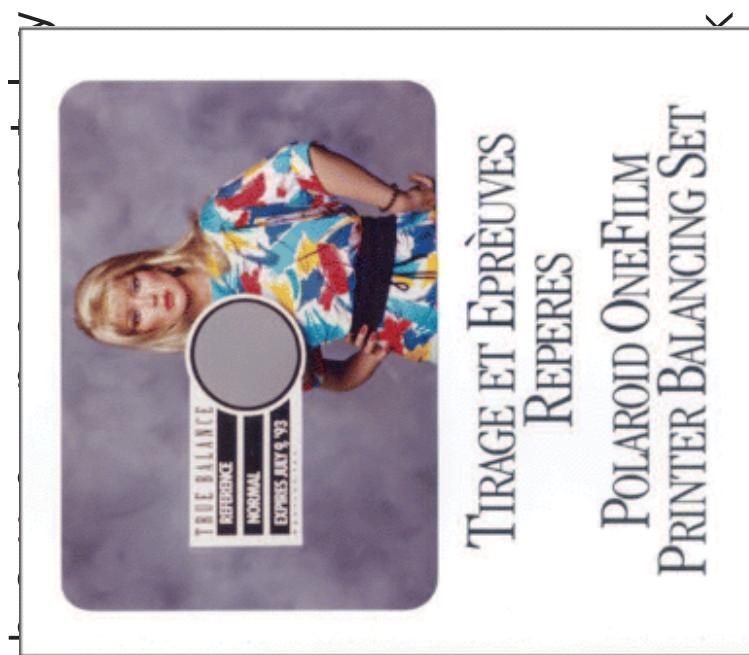
"Man is to computer programmer as () is to source code?"

- What sources of bias a people used in this exercise?
- Limited features
- Skewed sample
- Tainted examples
- Sample size disparity

Racials in photography

Smartphone manufacturers started advertising how accurate they photographed skin tones

This has been issued in photography, on demand at different situations.



TIRAGE ET EPREUVES
REPÈRES

POLAROID ONEFILM
PRINTER BALANCING SET

In pre-digital photography "Shirley cards" fo images. These cards were used in cameras. Contained Caucasians models in the 1990s.

Figure 1

Polaroid Shirley card (Printed with permission of Polaroid)

Source: Lorin²³a Roth

Racial bias in photography

This issue has also been digitized photography tone scale is largely uniform.

TABLE 1 Fitzpatrick Classification of Skin Types I through VI

| Type I | Type II | Type III | Type IV | Type V | Type VI |
|---|--|---|---|---|---|
| White skin. Always burns, never tans. | Fair skin. Always burns, tans with difficulty. | Average skin color. Sometimes mild burn, tan about average. | Light-brown skin. Rarely burns. Tans easily. | Brown skin. Never burns. Tans very easily. | Black skin. Heavily pigmented. Never burns, tans very easily. |

The continued use of tone scale is the main evidence in processing off images compared to the Fitzpatrick scale. Dr Eli Issakshvili mentioned in the Monkia study with evidence improved images in the Fitzpatrick scale.

Sourcer: ²⁵ We take

Racial bias in photography

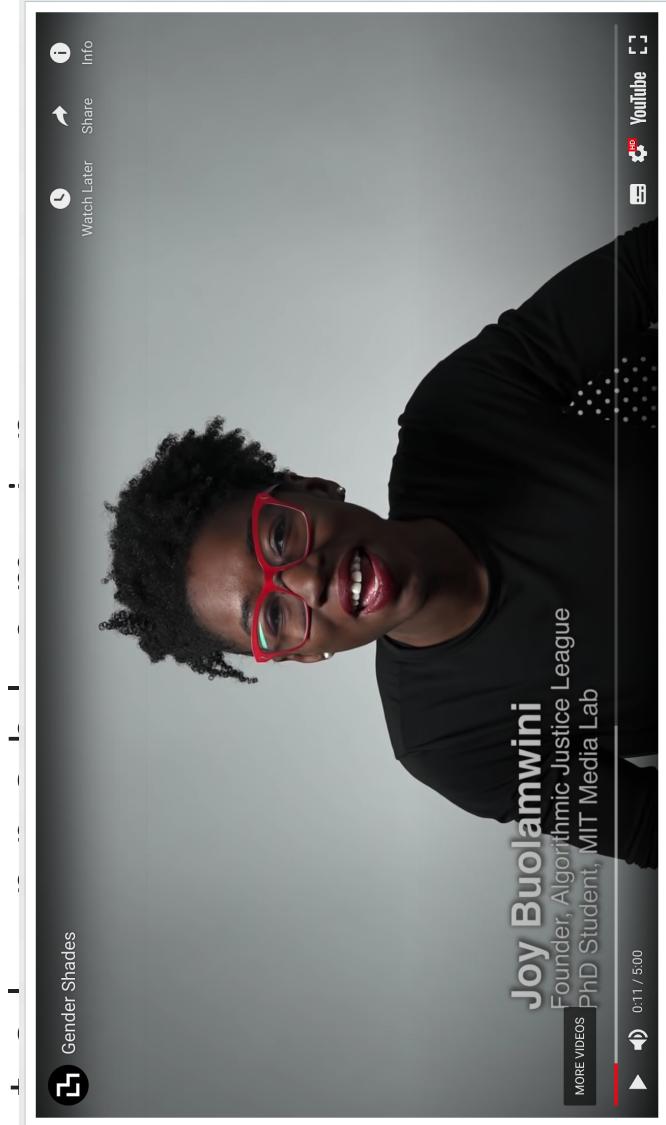
I'd like to recommend some additional resources:

- Vox [Messed it up](#) and [some yaws](#) built for [whi](#) atte i pte odpi sk ²⁶n.
- The Racial Bias Built [littles](#) Photography by [Lorraine Rosta](#) is a [raccoon](#) [page](#) [that](#) [is](#) [telly](#) [matrix](#) Norm: [image](#) [ecology](#) [ogies](#), [and](#) [cognition](#) [utive](#) [Equity](#).
- Google's [Resistant](#) [in](#) [Monk](#).

Racials in biass

This is a fundamental
ali face analysis
software.

Joy Buolamwini
better job of exposing
minutes than I car



Sourcendress.org

Racial bias in face photography

- What sources of bias are peaked in this study?
- Limited features
 - Skewed samples
 - Tailored expressiveness
 - Sample size differences

A C C O U N T a b n i d h s t p y a r & e r

Accountability & Transparency

These two terms are used in the analysis of how these two concepts fit together. I take two

These are liability issues²⁹ and transparency³⁰

The primary role of the organization is doing, it is not known what is regulated.

...

An important aspect is transparency and accountability and ethical obligations, accept responsibility transparent manner

... but the rest of the document is very heavy

Accontability

For a more thorough explanation of accountability³⁰

Transparency and game

Ist' Worthwhile mentioning that transparency algorithms is applied:

Also, in some cases, transparency may lie a system. "For example, even the minimal open faucet has a low level of certain topics by boots and coordinates near the end of transparency
Source: Cap 30 et al 2018

- How else could we game algorithms?

O t h e r s o u r c e s o f

Other sources of uncertainty

Earlier we tried to calculate the cost of one software bias:

- Proxies
- Limited features
- Skewed sample
- Attended examples
- Sample size disparities

These are almost always concerned with the training data behind algorithms.

These whole universe of ways we can bias an algorithm could

Fairness and Abstraction Systems

Selbst et al publish a paper on abstraction in 2008³¹ "abstraction for context for fairness".

These traps are designed to prevent a conflict for the interests between the technical system behind our algorithms and the social world in which they're applied.

Sell boost Abstract interpretation

Failure model the entire system over which be enforced

Example:

In the US criminal justice pipeline a defendant and perpetrator < br />

However - usually this is the risk occasionally considered

These risk assessment tools are presented on not account for³¹) could end up saying reetion b e t we

The intended application of the algorithm into account their use by judges and the measured

Sell boost Abstract Impact

Failure to understand how repurposing all goes in a context may be misleading, in a different context

Here's a really nice quote

Within computer science, it is considered to be used effectively in social contexts. "But what applies in Kentuckey really matters." Or you can't have others miss them if they're applied to employment. How we think about a failure

Source: K³²ren Hao

Sell best Abstract or: a p

Failure to account for the full meaning of be procedural, contextual, and contestable mathematical formalisms

I did promise we would not groer in other moral recruiting procedures than particular parties (and dates fit) because it's about him.

Fairness and discernment in philosophy have long debated. They are at the moment a complex concept that has been contested, and each party has its own definition. The concepts of fairness and proprieties are core elements of the moral system.

Sellbost Approach:

Failure to understand how the insertion of changes the behaviors and embedded values

Selbost³¹ identified a central focus of the intent of the justiciability of rights increases the intention of the law.

Sell best Abstraction: Trap

Failure to recognize the possibility that
technology

Modeling requires intentionality
politicality contended, moves
how it moves

C O M R S : A l g o r i t h m i
a s s e s s m e n t s

COMPARISON : Artificial intelligence assessment

The algorithmic risk assessments were discussed at Si. Scailed between situ~~ated~~ y

The algorithmic risk assessments were discussed at Si. Scailed between situ~~ated~~ y

Proposed by ³³ in 2016 of a test invented for the first time. A breakthrough authors

ACCURACY
RISK score: 65.2%
HUMAN*: 67.0%

FALSE POSITIVE*

Black: 37.1%
White: 27.2%

FALSE NEGATIVE*

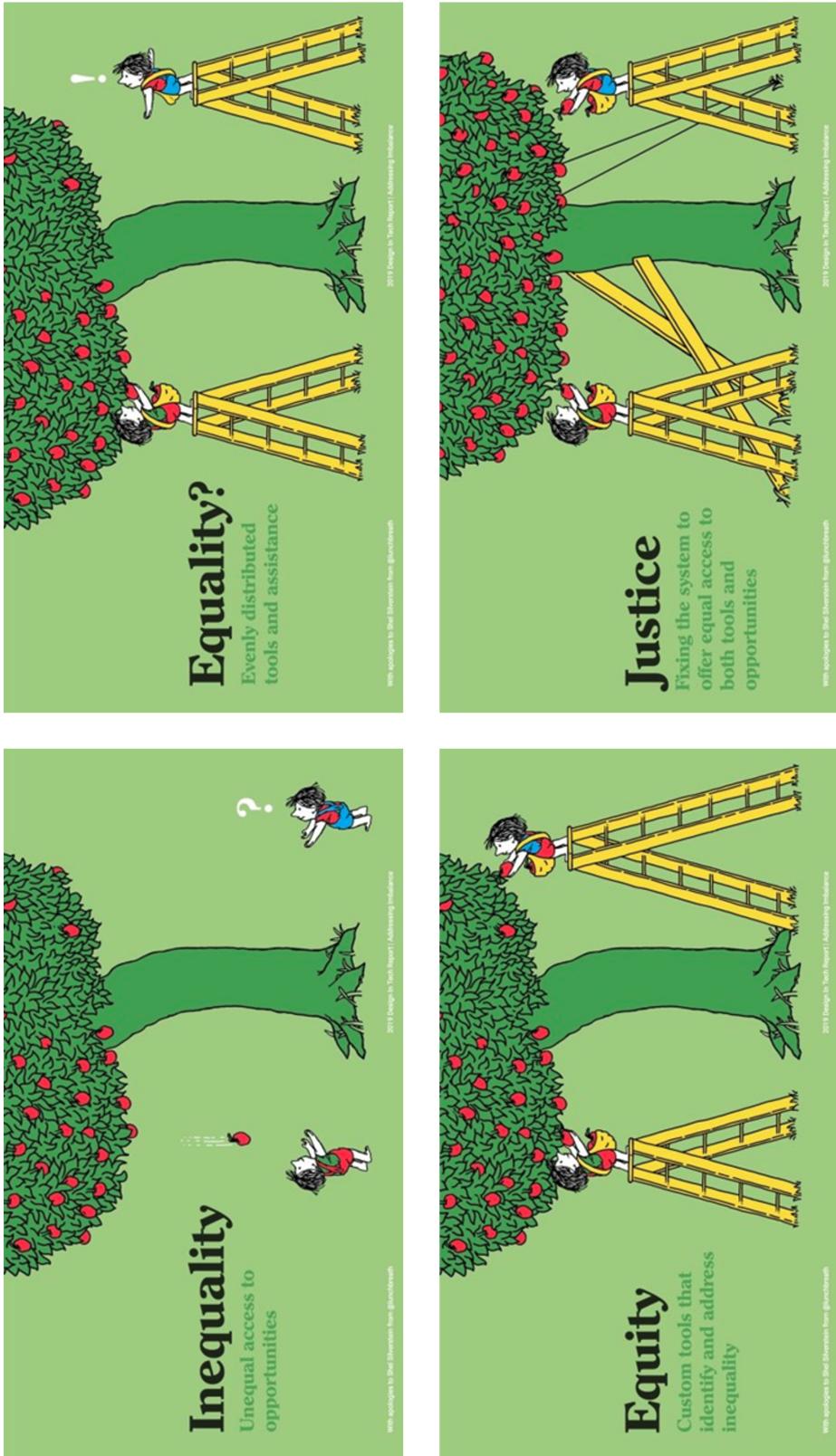
Black: 29.2%
White: 40.2%

The algorithm performance same and on Me cahrakin i ³⁵sael r ST
like Hanys hawt Fey D'ut
on th³⁶s topi c

Introducción a la teoría de los juegos

Inequality vs Justice

I thought it was important to not skip over these cartoony ~~Ma~~ ~~not~~ ~~re~~ ~~ne~~ ~~n~~ ~~t~~ ~~Dae~~ ~~tsi~~ ~~i~~ ~~ogn~~ ~~c~~ ~~on~~ ~~R~~ ~~E~~ ~~p~~ ~~o~~³⁷ ~~r~~ ~~t~~ ~~a~~ ~~n~~ ~~2~~ ~~O~~ designs by ~~R~~ ~~u~~ ~~t~~ ~~h~~



We have
that

stop

We have to stop some'

We've covered more than enough for your ass
complexities of ethics in algorithmic.
We didn't really discuss data privacy very
far information practices.

References

1. Suchol, I., & Schoni, I. M. "Less Than One"-Shot Learning from M < N Samples & Proceeding of the AAAI Conference on Artificial Intelligence 35, 9739–9746 (2021).
2. Uni et National General Assembly. UNI Assembly Declares the Right of Human Rights. (1948).
3. Awad, E., et al. **The Moral Machine Experiment**. Nature 563, 59–64 (2018).
4. Maxmen, A. **Self-driving cars do not understand moral choices**. Nature 562, 469–470 (2018).
5. Schoettl, B., & Sik, M. A Preliminary Analysis of Real World Crashes Involving Self-Driving Vehicles. National Transport Safety Board. Collected Between Vehicle Control by Developmental Automated Driving System and Pedestrian in Tempe, Arizona, March 18, 2018. (2018).
6. Hardt, M. Occupy AI白皮书: Who Are AI Systems Serving the 99%? (2013).
7. Barocas, S., & Selbst, A. D. Bias in Data Structures Impact (2016) doi [10.2139/ssrn.2477899](https://doi.org/10.2139/ssrn.2477899).
8. Landeau, A. Exploring Biometric Data. (2020).
9. Golde, J.C. Sources of Unintended Bias in Big Data. Mediamedia (2020).
10. Jerni, C., & Miltzou, B. F. T. Gaydar: Facebook finds itself exposed sexual orientation bias on Monday (2009) doi [10.5210/fm.v14i1.2611](https://doi.org/10.5210/fm.v14i1.2611).
11. Ensi, P. D., Frias-Martinez, V., Nevi, L., Sánchez-Beguer, C., & Venkatasubramanian, R. S. Runaway Feedback Loops in Predictive Policing (2017) doi [10.48550/arXiv.1706.09847](https://doi.org/10.48550/arXiv.1706.09847).
12. Lazer, D., Kennedy, R., Kiernan, G., & Vespi, Vanni, A. The Parable of Google Flu Trends in Big Data Analysis. Social Science Research 343, 1203–1205 (2014).
13. Giørg, J., et al. Detecting the influenza epidemic using search engine query data. Nature 457, 1012–1014 (2009).
14. Cook, S., Conrad, C., Fowlkes, A. L., & Mohebbi, M. H. Assessing Google Flu Trends Performance in the United States during the 2009 Influenza Virus A (H1N1) Pandemic. PLOS ONE 6, e23610 (2011).
15. Fraj, J., & László, Á. A literature review of Artificial Intelligence Impact on the Recruitment Process. International Journal of Engineering Management Sciences 6, 108–119 (2021).
16. Si, Q., & Mi, Y. The Experience and Nature of Racism Against Asian Supermarket Rations. Journal of Applied Psychology 90, 586–591 (2005).
17. Stauffer, J. M., & Buckley, M. R. The Experience and Nature of Racism Against Asian Supermarket Rations. Journal of Applied Psychology 90, 586–591 (2005).
18. Rogers, T., & Hoods, Programmers Believe About Names - What's in a Name? Shi, S. Journal of Applied Psychology 90, 586–591 (2005).
19. Rogers, T., & Hoods, Programmers Believe About Names - What's in a Name? Shi, S. Journal of Applied Psychology 90, 586–591 (2005).

20. Bol kbasi T., Chang, K.-W., Zou, J., Sal ragha, V. & Kal iaA. Man i \$ Computer Programmer as Woman i \$ Homemaker? Debi \$ Word Embedding. (2016) doi [10.48550/arXiv.1607.06520](https://doi.org/10.48550/arXiv.1607.06520).
21. Mi d@ T., Chen, K., Corrado, G. S. & Dean, J. A. Computer numeri representati nos of words i a hi lg-di kensi mal space. (2015).
22. Google .Real Tone on Google Pi at. Google Store (2022).
23. Roth, L. **Looki gnat Shi l ry the UI i tate Norm: Col w Bal race, Image Technol g\$ s and Cogni i t@E qui yt**. Canadian Journal of Communi ati 34, 111–136 (2009).
24. Monk, E. P., Jr. **The Uneasai gnificance of Col w ra: Ski tone Stratifica ti wi the Uni etd States**. Daedal 150, 76–90 (2021).
25. Ward, W. H., Lambreton, F., Goel N., Yu, J. Q. & Farma, J. M. TABLE 1, Fi zpatri kcci skificati wofSki types I through VI. (2017).
26. Vox. Col rofil mas bui fortwhi etpeopl .Here's whati di \$ dark ski .(2015).
27. Lewi ,S. The Raci l@Bi sBui Int\$ Photography. The New York Ti es (2019).
28. MIT Medi lab. Gender Shades. (2018).
29. European Parl ment Diectorate General for Parl ietary Research Servi es. A governance framework for al gti hmi c accountabi y and transparency. (Publ atic@ Office, 2019).
30. Capl @ R., Donovan, J., Hanson, L. & Matthews, J. AI gti hmi Accountabi y: A bri er. (2018).
31. Sel st, A. D., Boyd, D., Fri d@ reS. A., Venkatasubramani @ S. & Vertesi J. Fairness and Abstracti moi Soci techni at Systems. i Proceedi gs of the Conference on Fair ness, Accountabi y, and Transparency 59–68 (Associ @ wfor Computi gMachi ery, 2019). doi [10.1145/3287560.3287598](https://doi.org/10.1145/3287560.3287598).
32. Hao, K. Thi s howAI bi sareal happensand why i's so hard to fix. MIT Technol gy Revi w (2019).
33. Mattu, L. K., Jeff Larson. Machi enBi s ProPubl @ (2016).
34. Mattu, L. K., Jul iAgwi .How We Anal yzed the COMPAS Reci idivs Al gti hmi. ProPubl @ (2016).
35. Dressel J. & Far ,H. **The accuracy, fair ness, and li st of predi tc greci idivs**. Sci ike Advances 4, eaao5580 (2018).
36. TEDx Tal s The danger of predi tc eval gti hms i ori matj sjii ec| Hany Fari fTEDxAmoskeagMi larb. (2018).
37. Maeda, J. Presentati o Desi gni Tech Report 2019. Desi gni Tech (2019).