



<https://www.aclunc.org/blog/amazon-face-recognition-falsely-matched-28-members-congress-mugshots>

AI Fairness and Ethics

Lecturer: Dr. Xingjun Ma
School of Computer Science,
Fudan University
Fall, 2022

Recap: week 12

- Federated Learning
- Privacy in Federated Learning
- Robustness in Federated Learning
- Challenges and Future Research



Team/Project Registration



可信机器学习组队注册
扫一扫二维码打开或分享给好友



【腾讯文档】可信机器学习组队注册
<https://docs.qq.com/form/page/DU0dEYUp2R1BGa0NF>

It's time for team registration!



Biases in Current AI Models

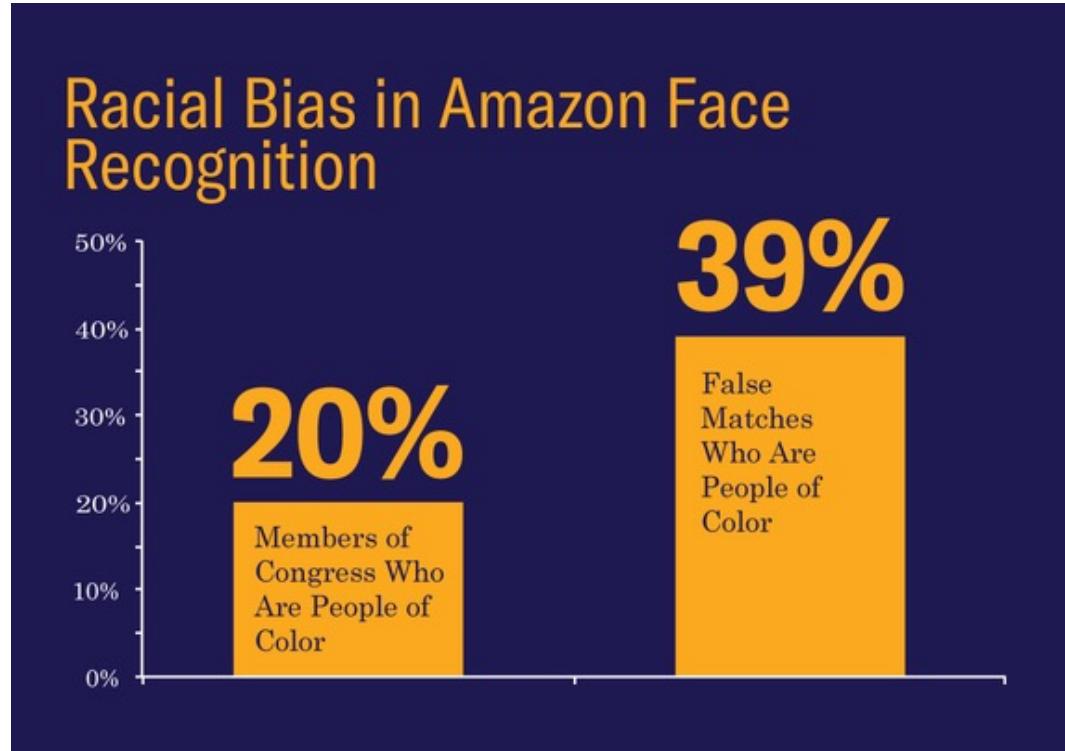


- A study conducted by ACLU (American Civil Liberties Union, 美国公民自由联盟)
- Object: Amazon facial recognition software Rekognition
- Methodology: the tool matches 28 Congress members with mugshots
- Mugshot database: 25,000 publicly available arrest photos
- The test costs only \$12.33 – less than a large pizza

<https://www.aclunc.org/blog/amazon-s-face-recognition-falsely-matched-28-members-congress-mugshots>



Biases in Current AI Models



- 20% of the members are people of color
- 39% of the matched criminals are people of color

<https://www.aclunc.org/blog/amazon-s-face-recognition-falsely-matched-28-members-congress-mugshots>



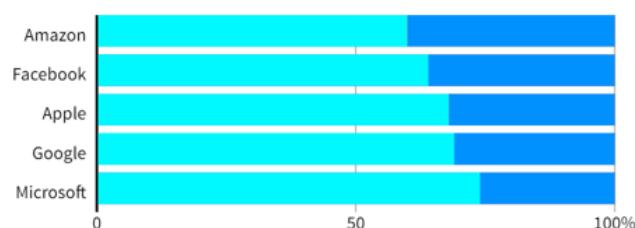
Biases in Current AI Models

Dominated by men

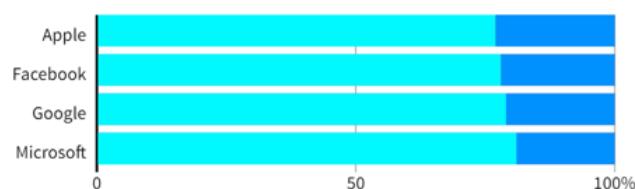
Top U.S. tech companies have yet to close the gender gap in hiring, a disparity most pronounced among technical staff such as software developers where men far outnumber women. Amazon's experimental recruiting engine followed the same pattern, learning to penalize resumes including the word "women's" until the company discovered the problem.

GLOBAL HEADCOUNT

Male Female



EMPLOYEES IN TECHNICAL ROLES



Note: Amazon does not disclose the gender breakdown of its technical workforce.
Source: Latest data available from the companies, since 2017.



- A machine learning based resume filtering tool
- It takes 100 resumes and returns the top-5
- It was then found to recommend only men for certain jobs

<https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>



Biases in Current AI Models



<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

- **COMPAS** (Correctional Offender Management Profiling for Alternative Sanctions) algorithm used by US court systems
- A fairness study conducted by ProPublica (a Pulitzer Prize-winning non-profit news organization)
- The prob of reoffend: black offenders (45%) vs white offenders (23%)

<https://towardsdatascience.com/real-life-examples-of-discriminating-artificial-intelligence-cae395a90070>



Biases in Current AI Models



Photo by M. Spencer Green / AP

arXiv > cs > arXiv:1706.09847

Computer Science > Computers and Society

[Submitted on 29 Jun 2017 (v1), last revised 22 Dec 2017 (this version, v2)]

Runaway Feedback Loops in Predictive Policing

Danielle Ensign, Sorelle A. Friedler, Scott Neville, Carlos Scheidegger, Suresh Venkatasubramanian

Predictive policing systems are increasingly used to determine how to allocate police across a city in order to best prevent crime. Discovered crime data (e.g., arrest counts) are used to help update the model, and the process is repeated. Such systems have been empirically shown to be susceptible to runaway feedback loops, where police are repeatedly sent back to the same neighborhoods regardless of the true crime rate. In response, we develop a mathematical model of predictive policing that proves why this feedback loop occurs, show empirically that this model exhibits such problems, and demonstrate how to change the inputs to a predictive policing system (in a black-box manner) so the runaway feedback loop does not occur, allowing the true crime rate to be learned. Our results are quantitative: we can establish a link (in our model) between the degree to which runaway feedback causes problems and the disparity in crime rates between areas. Moreover, we can also demonstrate the way in which (emph{reported}) incidents of crime (those reported by residents) and (emph{discovered}) incidents of crime (i.e., those directly observed by police officers dispatched as a result of the predictive policing algorithm) interact: in brief, while reported incidents can attenuate the degree of runaway feedback, they cannot entirely remove it without the interventions we suggest.

Comments: Extended version accepted to the 1st Conference on Fairness, Accountability and Transparency, 2018. Adds further treatment of reported as well as discovered incidents

Subjects: Computer and Society (cs.CY); Machine Learning (stat.ML)

Cite as: arXiv:1706.09847 [cs.CY]

- **PredPol** (predictive policing) algorithm biased against minorities.
- It predicts where crimes will occur in the future, designed to reduce human bias.
- It is already used by the USA police in California, Florida, Maryland, etc.
- It repeatedly sends police patrol to regions that contains a large number of racial minorities.

<https://towardsdatascience.com/real-life-examples-of-discriminating-artificial-intelligence-cae395a90070>



Biases in Current AI Models

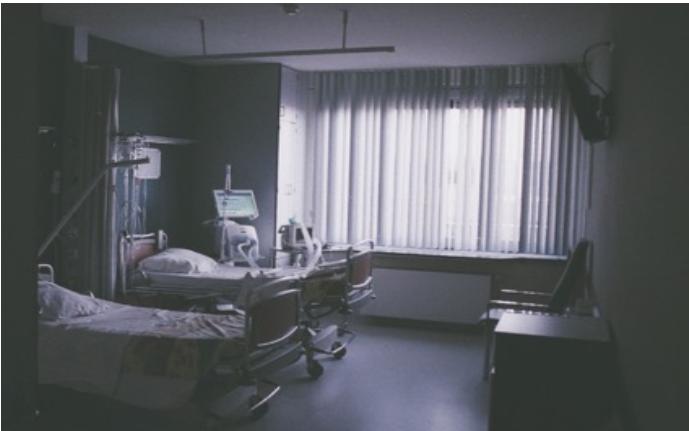


Photo by Daan Stevens on Unsplash

The screenshot shows the Science journal website. At the top, there's a navigation bar with links for NEWS, CAREERS, COMMENTARY, JOURNALS, and a search bar. Below the navigation is the Science logo. The main content area features a news article titled "Dissecting racial bias in an algorithm used to manage the health of populations". The article is categorized as a RESEARCH ARTICLE. Below the title, author information is listed: ZIAD OBERMEYER, BRIAN POWERS, CHRISTINE VOGELI, and SENDHIL MULLAINATHAN. The article was published on 25 Oct 2019, Vol 366, Issue 6464, pp. 447-453, DOI: 10.1126/science.aax2342. There are also social media sharing icons and a "CHECK ACCESS" button.

- Health care risk-prediction algorithm used for 200 million people in US hospitals predicts who needs extra health care.
- The algorithm heavily favours white patients over black patients, although race is not a variable for prediction.
- It was actually caused by a cost variable (black patients incurred lower health-care costs).

<https://www.scientificamerican.com/article/racial-bias-found-in-a-major-health-care-risk-algorithm/>



Biases in Current AI Models



眼睛太小被辅助驾驶系统识别为“开车睡觉”

16:00 ① 24.0 KBPS HD 5G 4G 5G



智驾分

了解智驾安全

90 / 100

守护里程: 814km ①

展开 ▾

待处理 ③

我的记录

使用辅助驾驶时, 频繁分神

阅读智驾须知

-1分

2022/07/24

使用辅助驾驶时, 频繁分神

查看正确操作

-1分

2022/06/26

使用辅助驾驶时, 长时间分神

查看正确操作

-2分

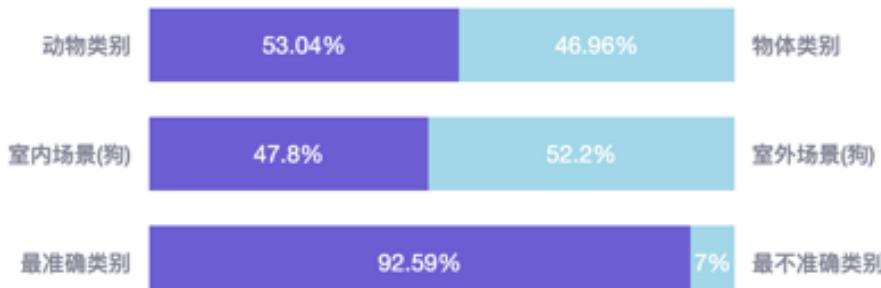
2022/06/26

<https://www.brookings.edu/research/enrollment-algorithms-are-contributing-to-the-crises-of-higher-education>

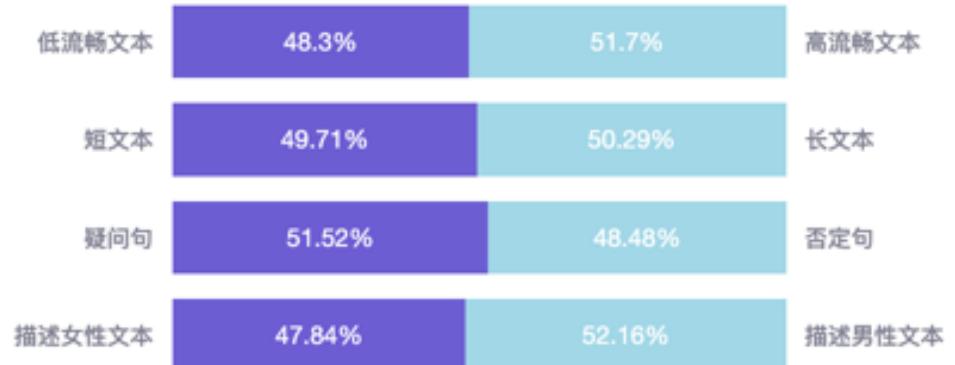


Demo展示

6 公平性



6 公平性



公平性 解释

类间不公平性: 选择两组 (类) 数据进行对比实验, 如果模型在这两组数据上预测置信度一致, 则结果是公平的, 否则结果差异越大, 公平性越差。**属性不公平性:** 选择两组 (同类不同属) 数据进行对比实验, 如果模型在这两组数据上预测结果一致, 则结果是公平的, 否则结果差异越大, 公平性越差。

图像分类模型

文本分类模型

<https://tech.openeglab.org.cn:8001/eval/image?code=Normal>



Definition of Bias

Fairness: the absence of any prejudice or favouritism toward an individual or group based on their inherent or acquired characteristics. (无差别化决策)

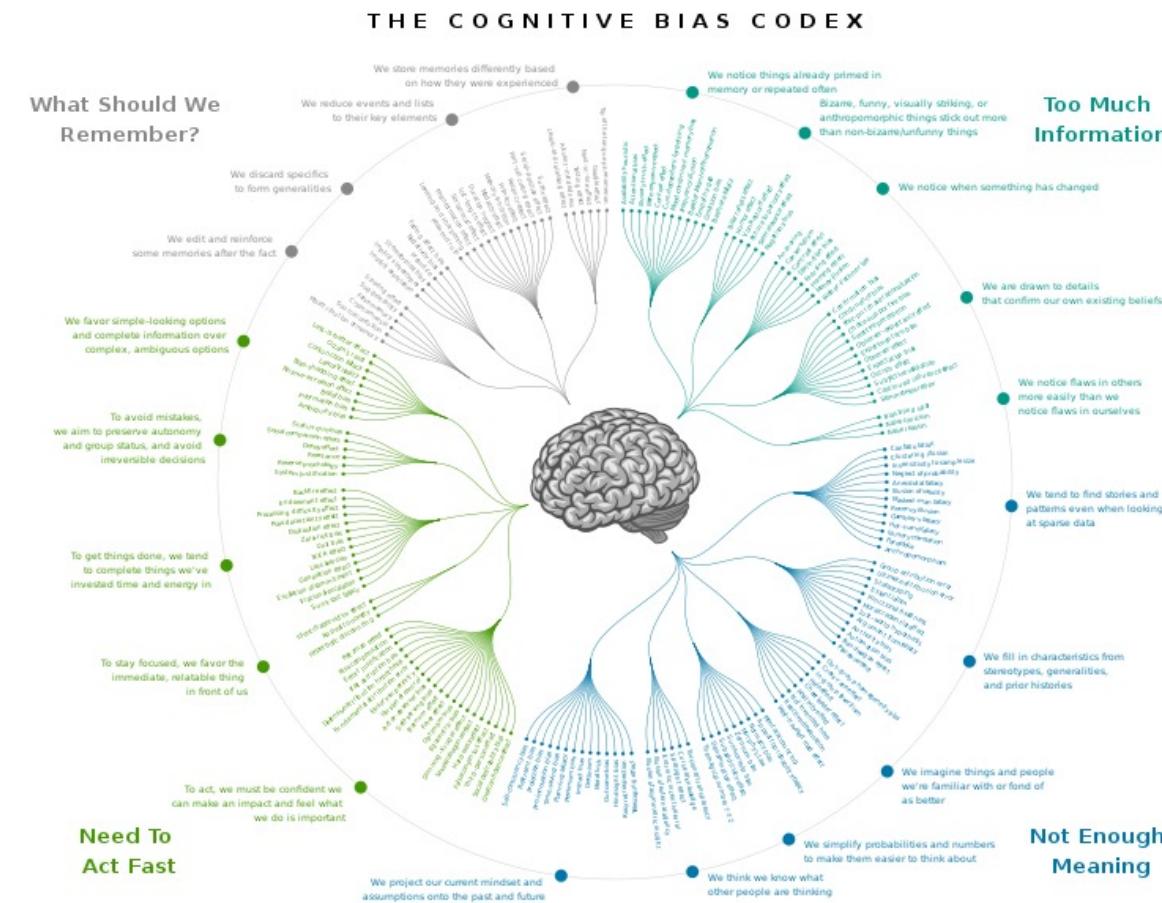
Bias: decisions are skewed toward a particular group of people not based on their inherent characteristics. (差别化决策)

Bias consists of attitudes, behaviors, and actions that are prejudiced in favor of or against one person or group compared to another. (社会学)

Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys (CSUR) 54.6 (2021): 1-35.
<https://diversity.nih.gov/sociocultural-factors/implicit-bias>



Cognitive Biases: Psychology and Sociology

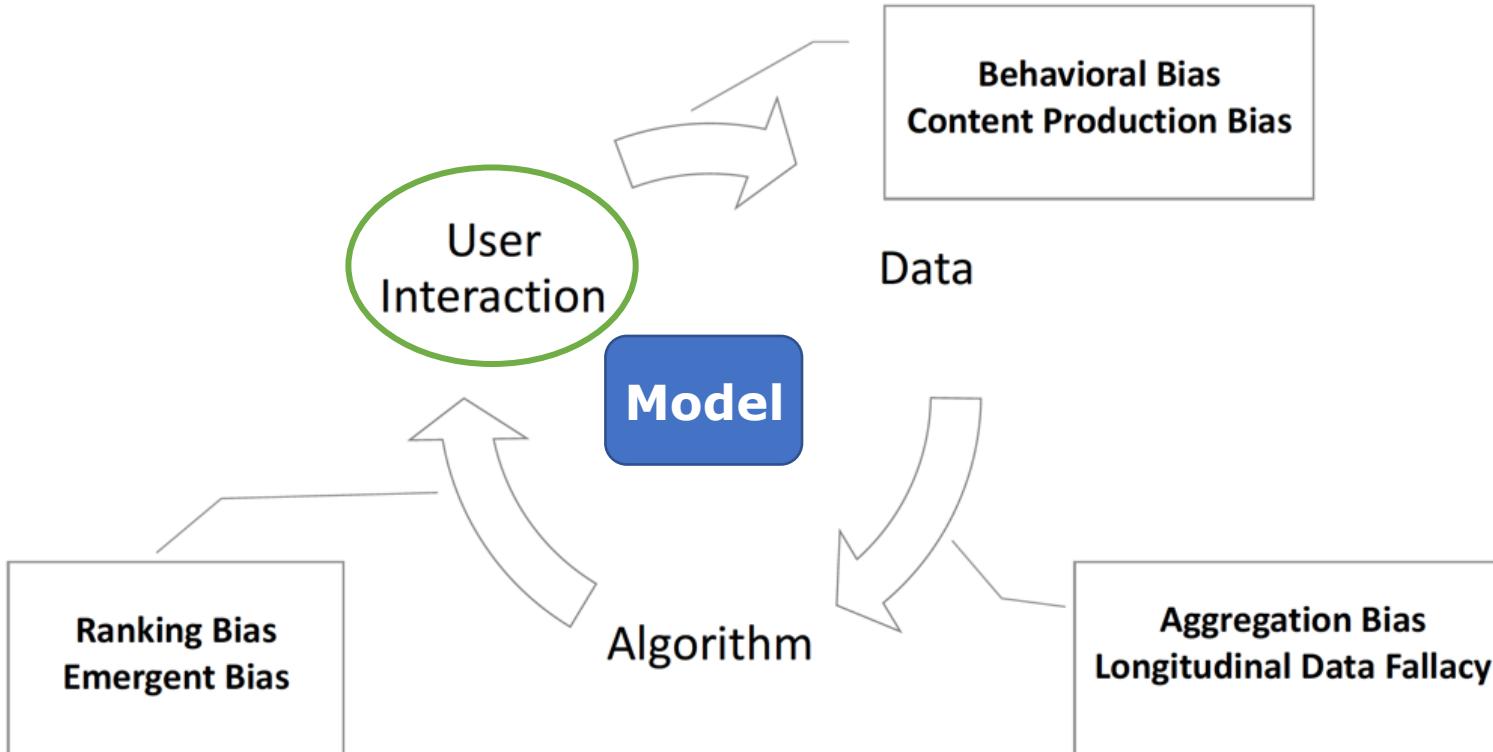


https://en.wikipedia.org/wiki/Cognitive_bias

https://upload.wikimedia.org/wikipedia/commons/6/65/Cognitive_bias_codex_en.svg



Types of Machine Learning Bias



[Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys \(CSUR\) 54.6 \(2021\): 1-35.](#)



Data Bias

- Measurement Bias
 - COMPAS 使用被捕次数和家庭成员被捕次数作为风险预测属性
- Omitted Variable Bias
 - 竞争对手（忽略的因素）的出现导致大量用户退订
- Representation Bias
 - 数据集的分布不具有全局代表性：比如ImageNet的地域分布
- Aggregation Bias
 - a) Simpson's Paradox : 在某个条件下的两组数据，分别讨论时都会满足某种性质，可是一旦合并考虑，却可能导致相反的结论
 - b) Modifiable Areal Unit Problem (MAUP) : 分析结果随基本面积单元（栅格细胞或粒度）定义的不同而变化的问题
- Sampling Bias : 跟representation bias类似，源自非随机采样
- Longitudinal Data Fallacy (纵向数据错误) : 未考虑时间因素
- Linking Bias : 社交网络图里面用户交互规律和连接关系有很大不同

[Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys \(CSUR\) 54.6 \(2021\): 1-35.](#)



Algorithmic Bias

- Algorithmic Bias

- 优化、正则化方法，统计分析方法，对数据的有偏使用

- Recommendation Bias

- 呈现方式和排行顺序存在偏见

- Popularity Bias

- 越流行的物体得到的推荐越多，进而获得更多的点击

- Emergent Bias :

- 软件完成设计后用户群体已经变了

- Evaluation Bias :

- 使用不恰当的基准数据集去衡量模型

[Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys \(CSUR\) 54.6 \(2021\): 1-35.](#)



User Bias

- Historical Bias
 - 历史数据存在偏见，比如搜索“女CEO”会根据历史数据返回很少的女性
- Population Bias
 - 平台用户群体不同，比如女生喜欢用Pinterest, Facebook, Instagram，而男生喜欢用Reddit or Twitter
- Self-Selection Bias
 - 采样偏见的一种，比如对于意见调查
- Social Bias：
 - 别人的行为影响我们的决定（别人都给高分，你给不给？）
- Behavioral Bias：
 - 不同圈子/平台上的人的行为不同，比如emoji表情的使用习惯
- Temporal Bias：
 - 人群和行为都会随时间而变化，比如twitter上有时候会用hashtag有时又不用
- Content Production Bias：
 - 每个人创造内容的方式和习惯不同，比如不同群体的文字使用习惯不同

[Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys \(CSUR\) 54.6 \(2021\): 1-35.](#)



Existing Bias Datasets

| Dataset Name | Size | Type | Area |
|--|-----------|----------------|------------------------|
| UCI adult dataset | 48,842 | income records | Social |
| German credit dataset | 1,000 | credit records | Financial |
| Pilot parliaments benchmark dataset | 1,270 | images | Facial Images |
| WinoBias | 3,160 | sentences | Coreference resolution |
| Communities and crime dataset | 1,994 | crime records | Social |
| COMPAS Dataset | 18,610 | crime records | Social |
| Recidivism in juvenile justice dataset | 4,753 | crime records | Social |
| Diversity in faces dataset | 1 million | images | Facial Images |
| CelebA | 162,770 | images | Facial Images |



公平性定义

Def. 1: Equalized Odds

- 同等机会对，同等机会错

Def. 2: Equal Opportunity

- 同等机会对

Def. 3: Demographic Parity

- 个体存在与否不影响对

Def. 4: Fairness Through Awareness

- 输入相近，结果相同

Def. 5: Fairness Through Unawareness

- 决策不适用偏见属性

Def. 6: Treatment Equality

- 错误的数量一直

[Mehrabi et al. "A survey on bias and fairness in machine learning." ACM Computing Surveys \(CSUR\) 54.6 \(2021\): 1-35.](#)



Fair Machine Learning: Dataset Description

Show dataset statistics and creation details



Dataset and creation
details

[Gebru, Timnit, et al. “Datasheets for datasets.” Communications of the ACM 64.12 \(2021\): 86-92.](#)



Fair Machine Learning: Dataset Labels

Use dataset specifications (数据集说明书)

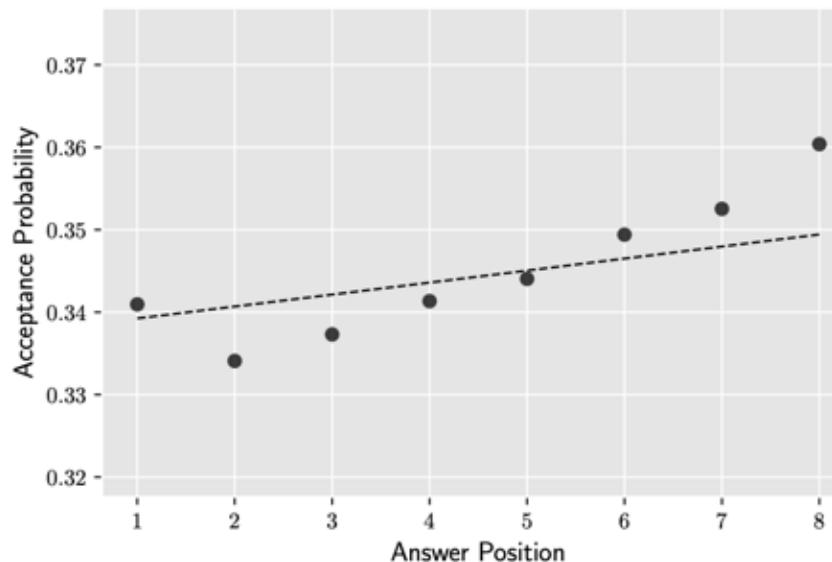


[Gebru, Timnit, et al. “Datasheets for datasets.” Communications of the ACM 64.12 \(2021\): 86-92.](#)

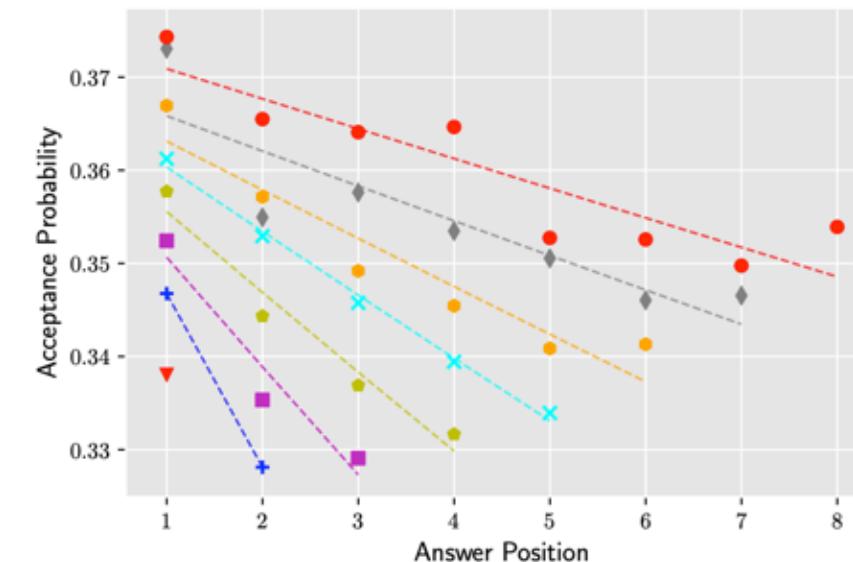


Fair Machine Learning: Dataset Labels

□ Simpson's paradox testing (合在一起结论变了)



(a) Aggregated Data



(b) Disaggregated Data

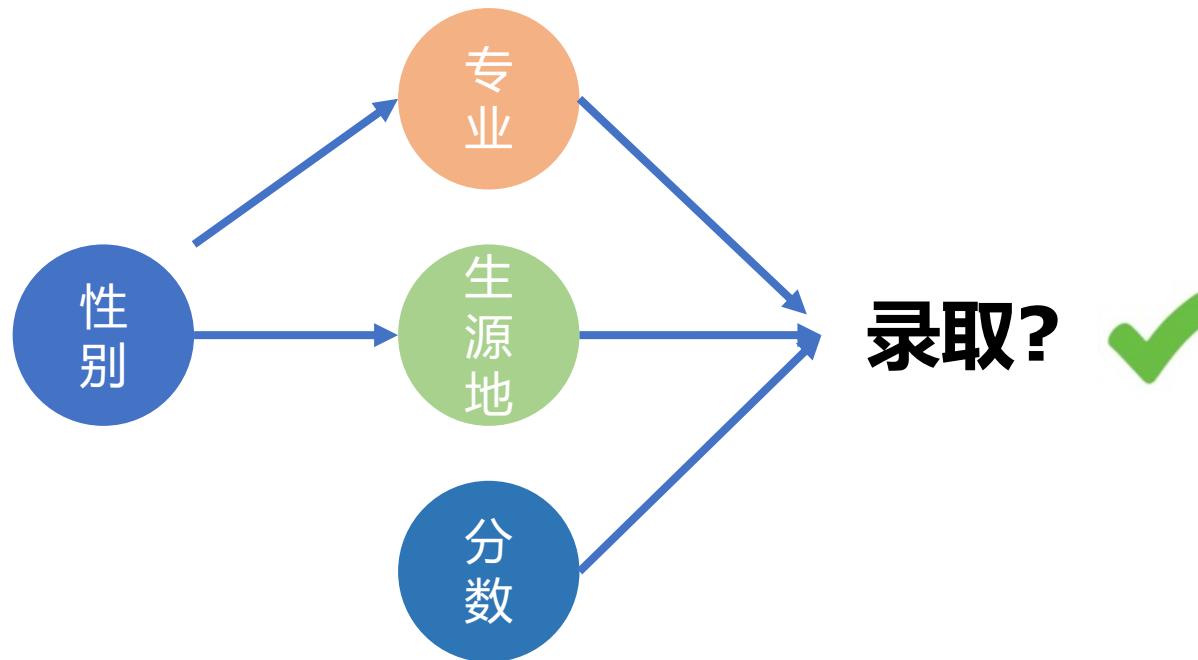
Stack Exchange: 第几个回答更容易被接受为最佳答案？ (b) : 基于session length划分group

[Gebru, Timnit, et al. "Datasheets for datasets." Communications of the ACM 64.12 \(2021\): 86-92.](#)



Fair Machine Learning: Causality

□ Identify and remove biases with causal graphs

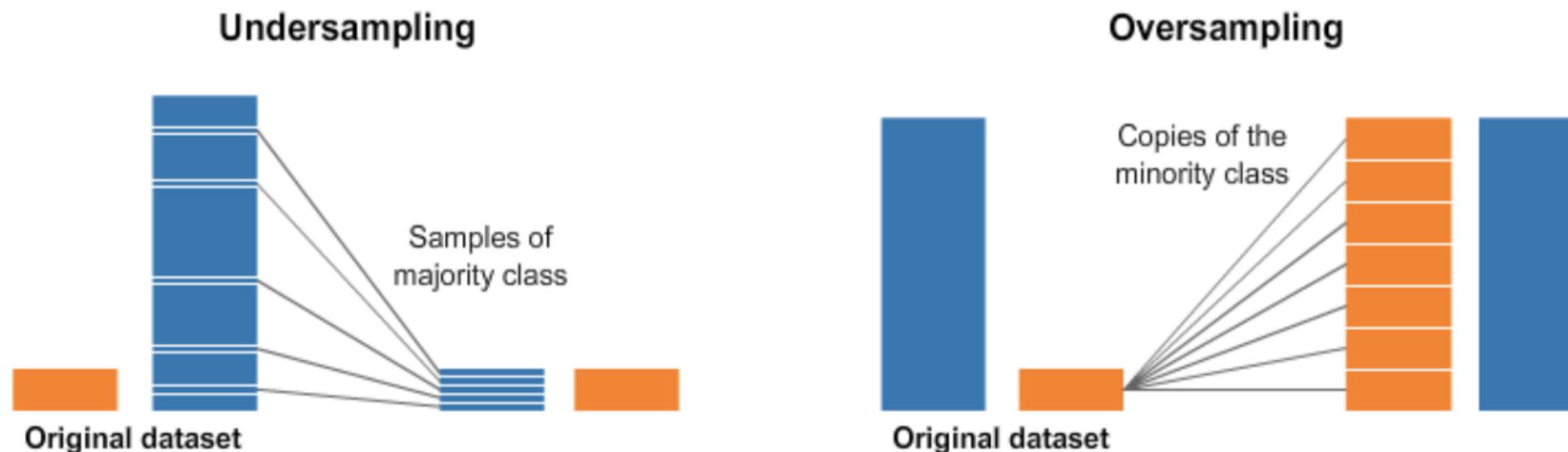


[Zhang, Lu, Yongkai Wu, and Xintao Wu. "Achieving non-discrimination in data release." SIGKDD, 2017.](#)



Fair Machine Learning: Sampling/Re-sampling

□ Balancing the minorities v majorities by re-sampling

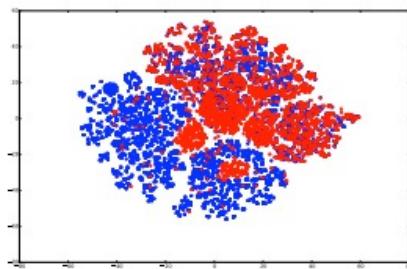


<https://www.kaggle.com/code/rafjaa/resampling-strategies-for-imbalanced-datasets/notebook>

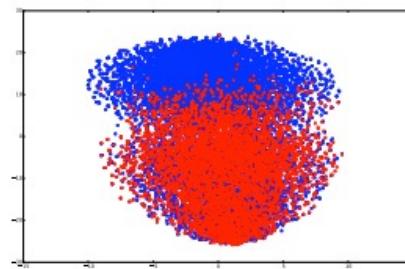


Fair Machine Learning: Fair Representations

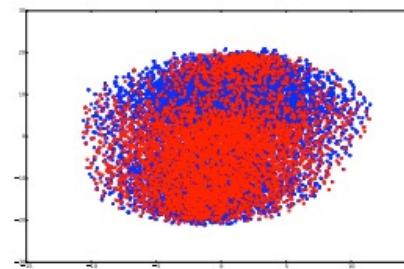
□ Fair AutoEncoder



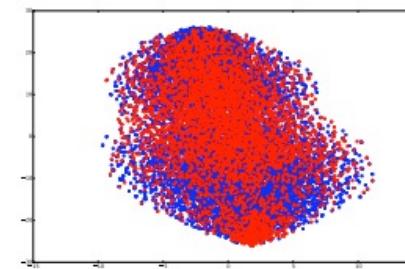
(a)



(b)



(c)



(d)

原始数据

MMD wo s

s wo MMD

s + MMD

Blue: male; **Red:** female

[Louizos, Christos, et al. "The variational fair autoencoder." arXiv preprint arXiv:1511.00830 \(2015\).](https://arxiv.org/abs/1511.00830)



Fair Machine Learning: debiasing

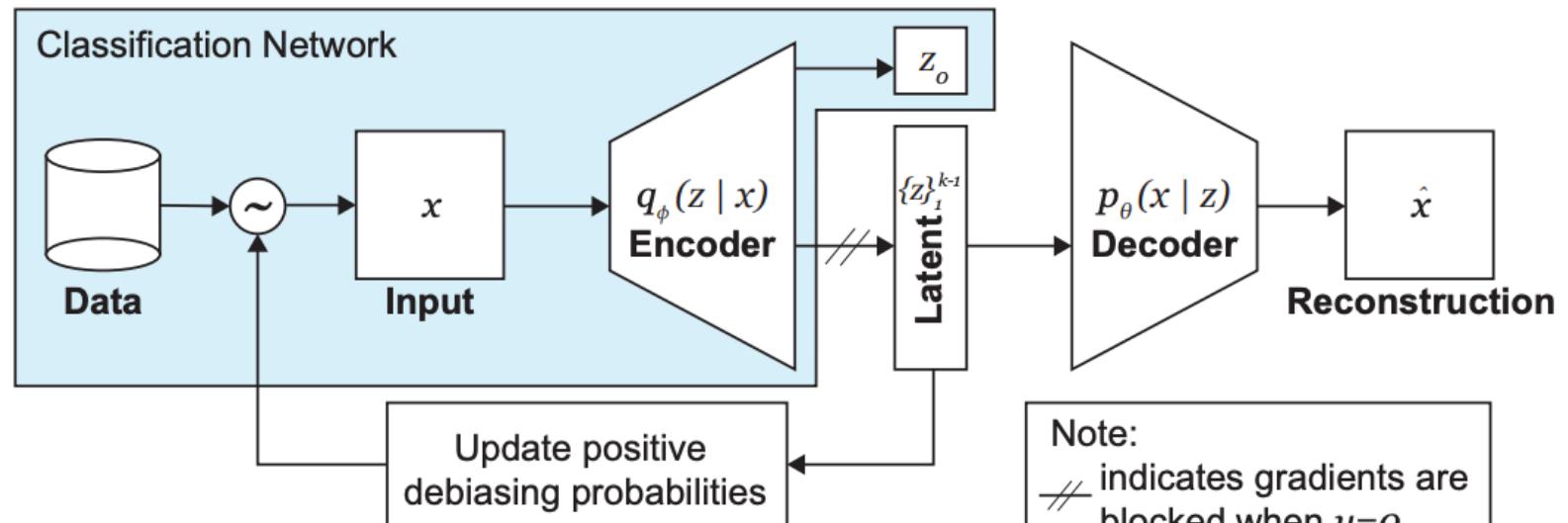
□ Debiasing Variational Autoencoder (DB-VAE)

Random Batch Sampling During Standard Face Detection Training

Homogenous skin color, pose
Mean Sample Prob: 7.57×10^{-6}

Batch Sampling During Training with Learned Debiasing

Diverse skin color, pose, illumination
Mean Sample Prob: 1.03×10^{-4}

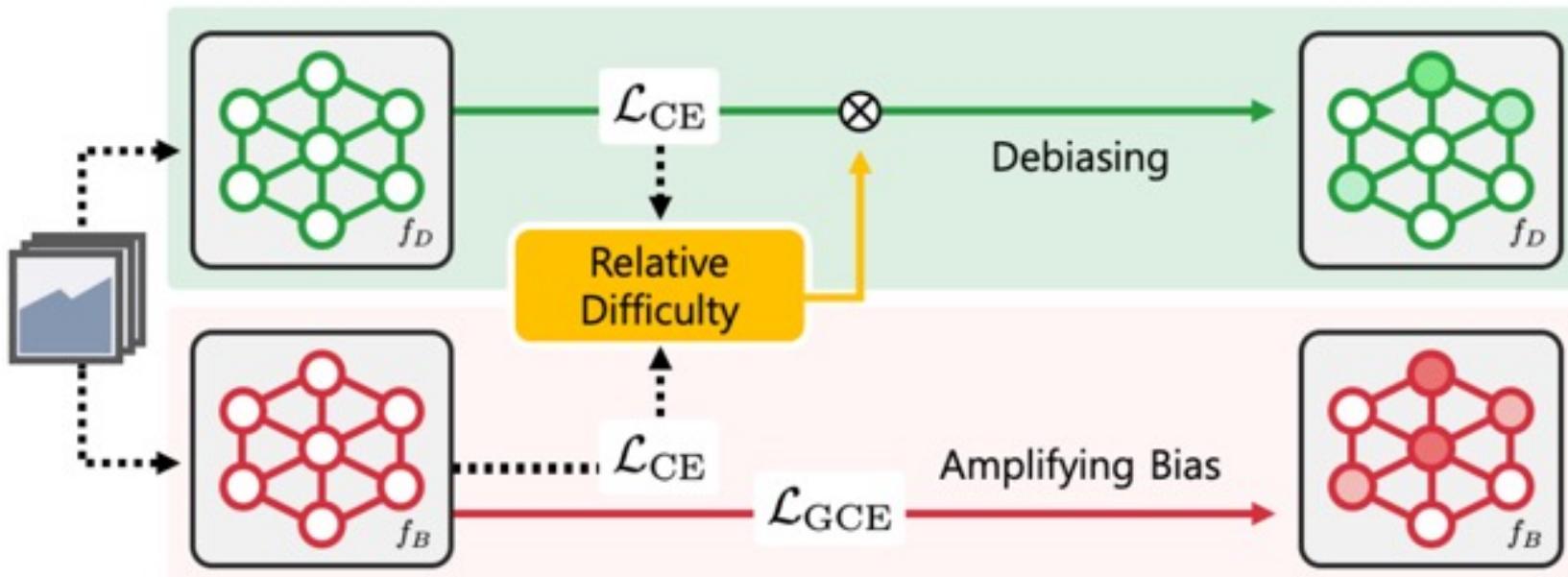


[Amini, Alexander, et al. "Uncovering and mitigating algorithmic bias through learned latent structure." AAAI. 2019.](#)



Fair Machine Learning: Unbiased Learning

□ De-biasing classifier from biased classifier

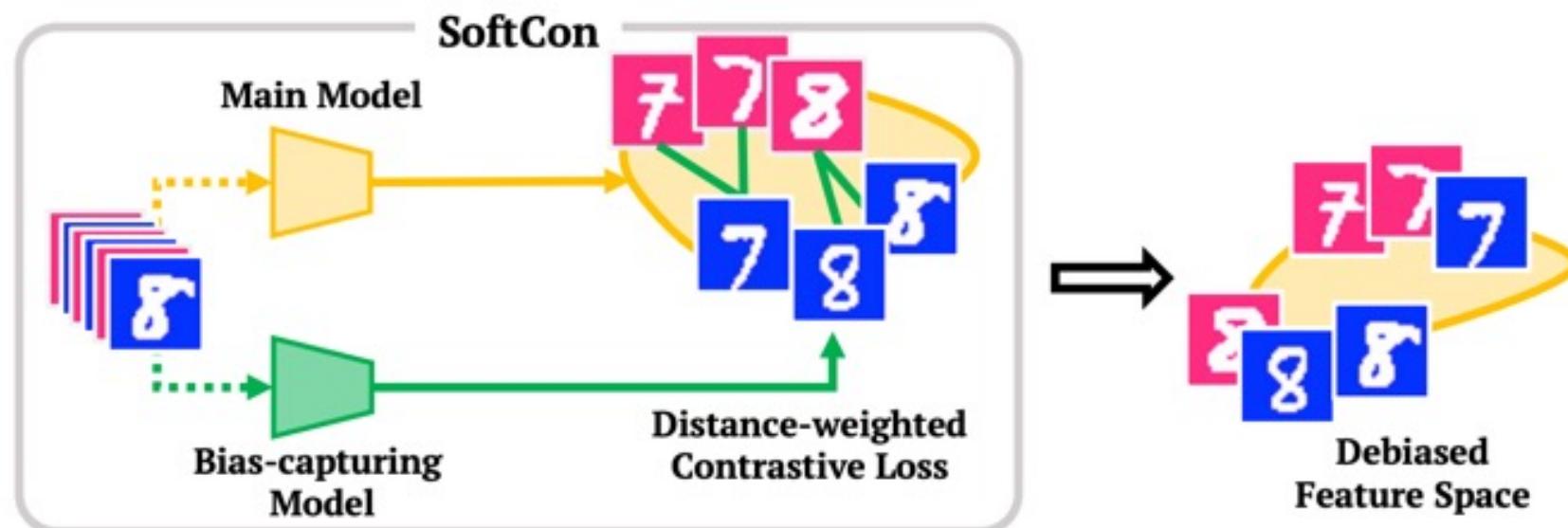


[Nam, Junhyun, et al. "Learning from failure: De-biasing classifier from biased classifier." NeurIPS, 2020.](#)



Fair Machine Learning: Unbiased Learning

□ Unbiased classification with bias capturing model

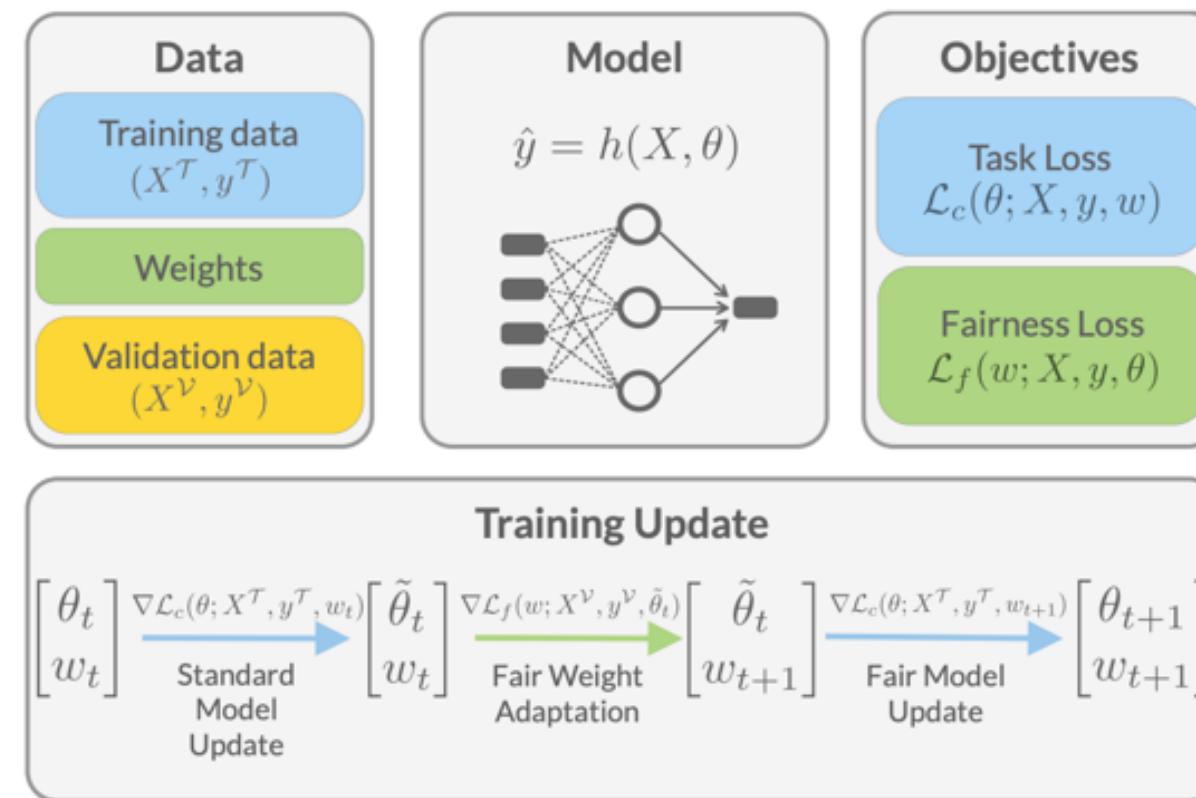


[Hong, Youngkyu, and Eunho Yang. "Unbiased classification through bias-contrastive and bias-balanced learning."](#) NeurIPS, 2021.



Fair Machine Learning: Reweighting

□ FORML: Learning to Reweight Data for Fairness

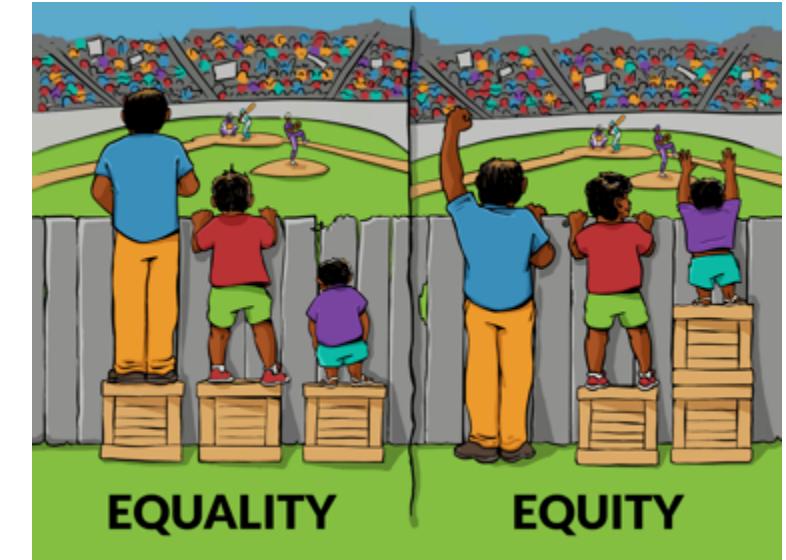


<https://machinelearning.apple.com/research/learning-to-reweight-data>



Remaining Challenges

- **Bias mining**: how to automatically identify biases for a given dataset and model
- **A general definition of bias/fairness**: a ML and societal definition of bias
- **From equality to equity**: 从平等到公平
- **Efficient fair learning**: fine-tuning based
- **In-situ debiasing**: identify and fix bias on-site





伦理规范、道德规范、职业道德



**“Ethics are moral principles that govern a person's or group's behaviour.
Synonyms: moral code, morals, morality, values, rights and wrongs, principles,
ideals, standards (of behaviour), value systems, virtues, dictates of conscience”**

Source: Oxford Dictionaries



Ethics, morals and rights - definitions

- **Ethics** – the study of the general nature of morals and of the specific moral choices to be made by the individual in his/her relationship with others. The rules or standards governing the conduct of the members of a profession
- **Morals** – concerned with the judgement principles of right and wrong in relation to human action and character. Teaching or exhibiting goodness or correctness of character and behaviour.
- **Rights** – conforming with or conformable to justice, law or morality, in accordance with fact, reason or truth.



Ethics – multiple meanings

- Philosophical ethics : attempts to use reason in order to answer various kinds of ethical questions.
- Describes a particular person's own, individual principles or habits. For example: "John has good ethics."
- Characterizes the questions of right-conduct in some specific sphere, even when such right-conduct is not examined philosophically: "business ethics," or "the ethics of child-rearing" may refer, but need not refer, to a philosophical examination of such issues.
- Philosophical ethics, or "ethical theory," is not the exclusive use of the term "ethics" in English



Morals – confused ...

- *Ethics* (also known as *moral philosophy*) is the branch of philosophy that addresses questions of morality.
- The word 'ethics' is commonly used interchangeably with 'morality' ... and sometimes it is used more narrowly to mean the moral principles of a particular tradition, group, or individual.



Rights ... multiple understandings

... a right to life, a right to choose; a right to vote, to work, to strike; a right to one phone call, to dissolve parliament, to operate a forklift, to asylum, to equal treatment before the law, to feel proud of what one has done; a right to exist, to sentence an offender to death, to a distinct genetic identity; a right to believe one's own eyes, to pronounce the couple husband and wife, to be left alone, to go to hell in one's own way (Wikipedia)



Who teaches us what is ethical?



Who teaches us what is ethical?

- Holy Book
- Mama
- Preacher
- Teacher
- Lawyer
- Doctor
- Government

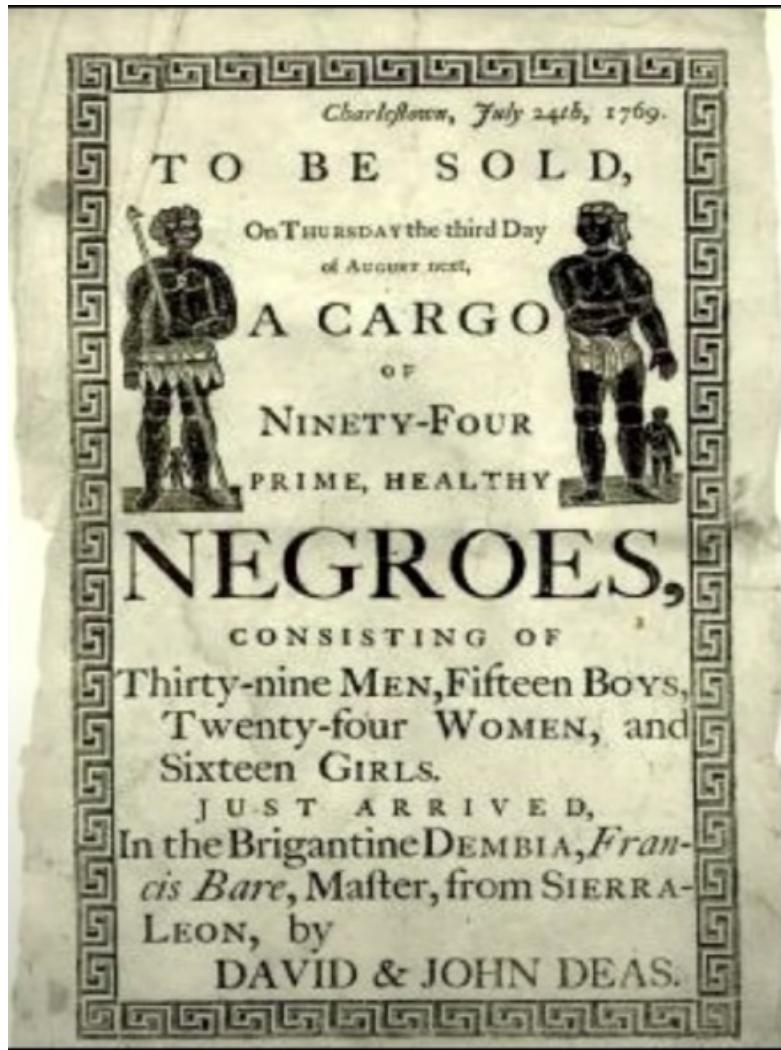




http://www.cityprofile.com/south-carolina/photos/12451-charleston-charleston_sc8.html

<https://www.youtube.com/watch?v=iiAirfn-lBI>





Bible



"Slaves, obey your earthly masters with fear and trembling"

Ephesians 6:5

"Tell slaves to be submissive to their masters and to give satisfaction in every respect"

Titus 2:9



Mama



THE
BLACK GAUNTLET:
A TALE OF
Plantation Life in South Carolina.

BY
MRS. HENRY R. SCHOOLCRAFT,
WIFE OF THE INDIAN CHIEF, AND AUTHOR OF "AFRICAN LETTERS," ETC. ETC.

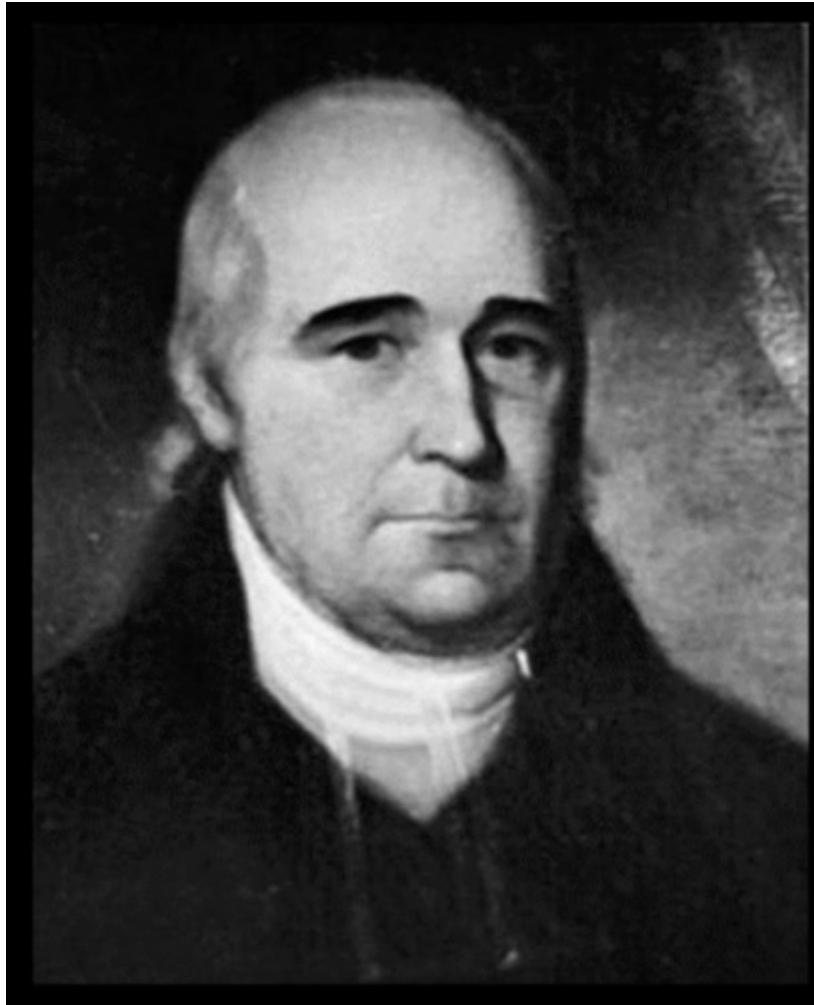
"God has placed a mark on the negro, as distinctive as that on Cain; and I do not believe there is a white man, woman, or child, on the face of the earth, who does not, in his deepest heart regard the African an inferior race to his own."

J. B. LIPPINCOTT & CO.
1860.

Mrs. Henry R. Schoolcraft, author of The Black Gauntlet: A Tale of Plantation Life in South Carolina
http://images.frick.org/PORTAL/IMAGEINFO.php?server=MTkyljE2OC4xMC43Mg==&siteurl=&file=/Volumes/Photoarchive/Duplicate_Negatives/21736_30104/21750_POST.tif



Preacher

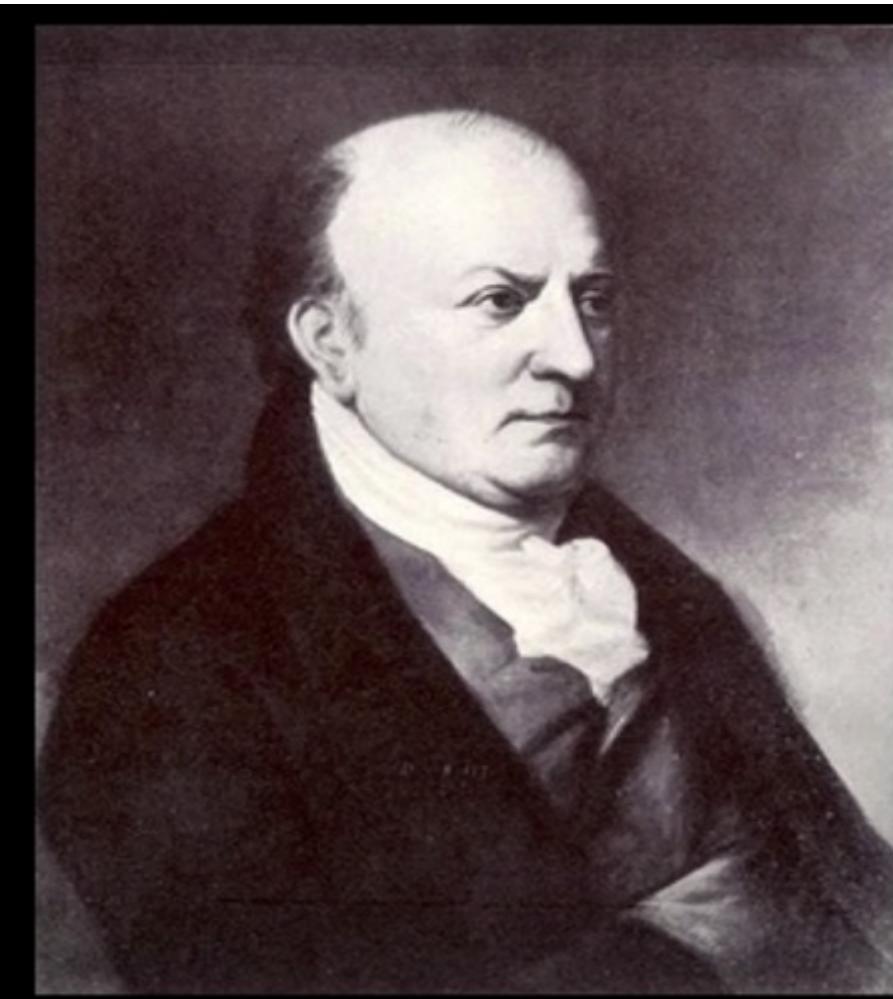


*The Most Reverend Richard Furman,
Baptist Pastor
Charleston, South Carolina.*

“The holding of slaves is justifiable by the doctrine and example contained in Holy writ; and is; therefore consistent with Christian uprightness...”



Teacher & Lawyer



*President Thomas Cooper,
University of South Carolina
Oxford student, lawyer,
philosopher, Chemistry professor*

1826 pamphlets

Outlined his belief that slave labor
was an economic necessity and that
the white race was superior.

http://library.sc.edu/digital/slaverysc/Presidents_Professors.html



Doctor

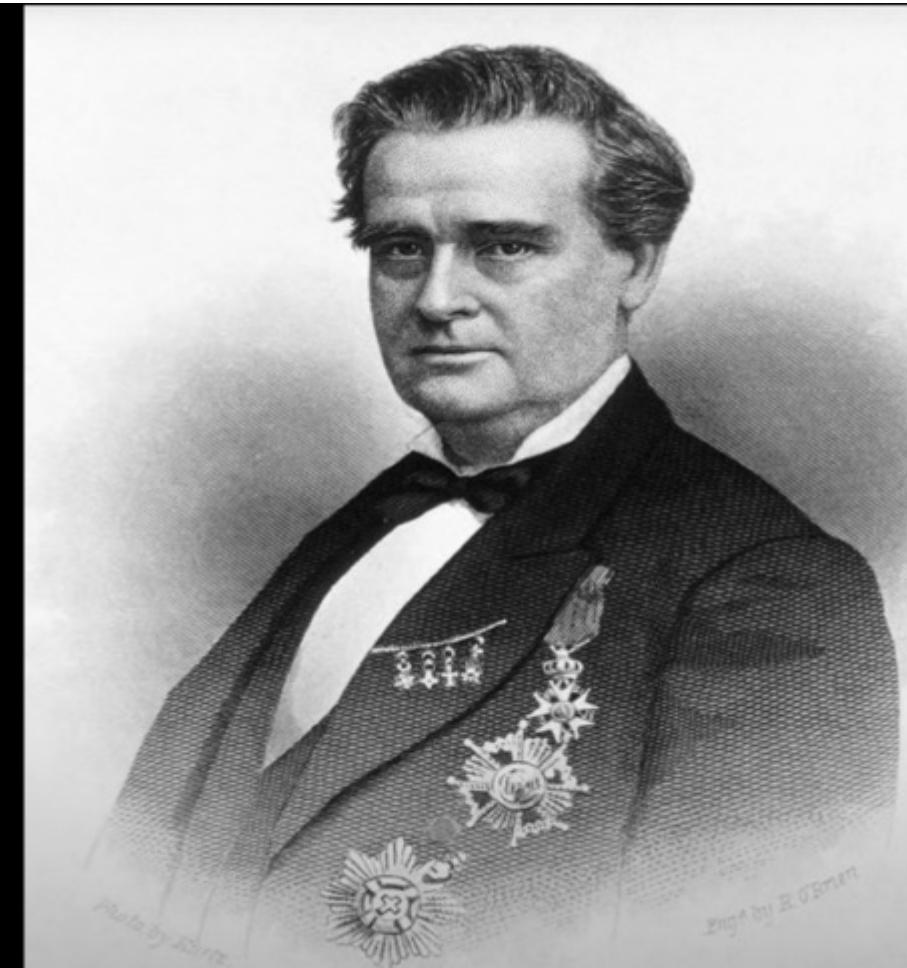
*Doctor J. Marion Sims
Founder of Gynecology*

No need for surgical anesthesia for
blacks or Irish...

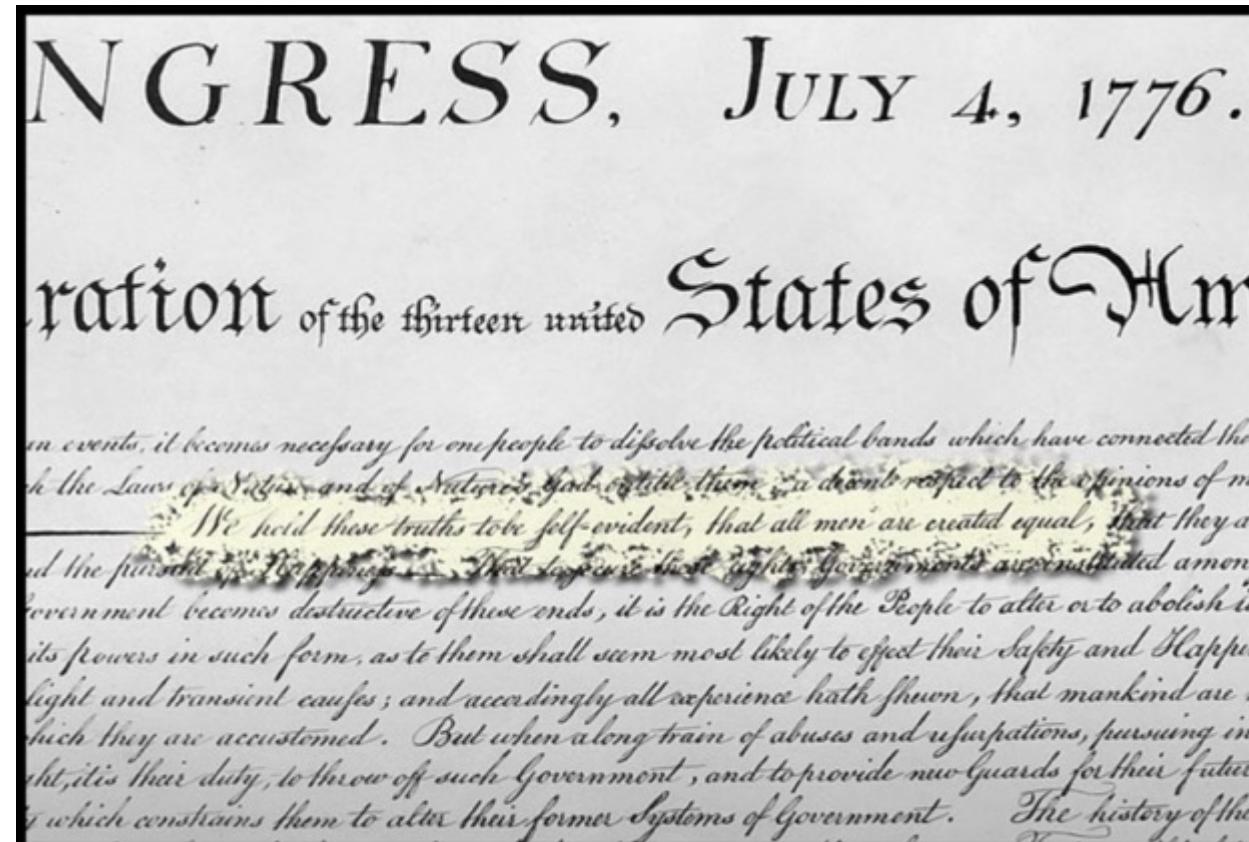
Bought Slaves to experiment on...

https://upload.wikimedia.org/wikipedia/commons/0/0b/James_Marion_Sims.jpg

<http://www.medscape.com/viewarticle/479892>



US Government



THE
SLAVERY CODE

OF THE
DISTRICT OF COLUMBIA,

WITH NOTES AND JUDICIAL DECISIONS EXPLANA-
TOGETHER
TORY OF THE SAME.

BY A MEMBER OF THE WASHINGTON BAR.

LC
WASHINGTON:
L. TOWERS & CO., PRINTERS.
1862.

Chuck Coker Flickr
Library of Congress <http://www.loc.gov/search/?q=Slave+Code&sp=3>



Today's ethics principles are built
upon the progress of human
civilization.



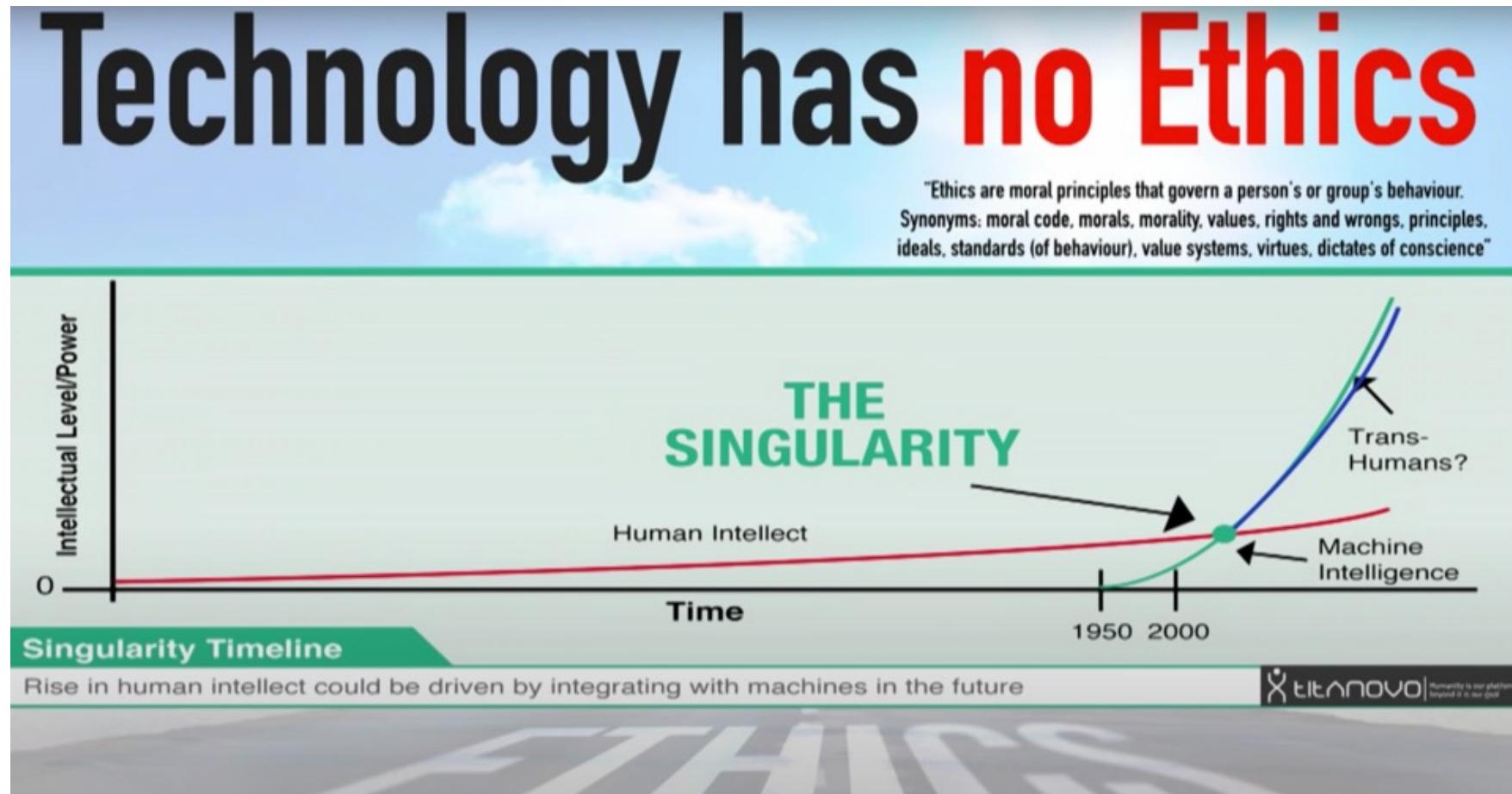
Ethics, technology and the future humanity?

AI Ethics

Laws and ethics are falling far behind modern technologies.



Unfortunately



<https://www.youtube.com/watch?v=bZn0IfOb61U>



FACT



"Don't be evil"

Google's motto, was suggested by Paul Buchheit.

WebDevelopersNotes.com



Google Removes ‘Don’t Be Evil’ Clause From Its Code Of Conduct

Share    

Kate Conger

Published 2 years ago: May 19, 2018 at 8:00 am - Filed to: ALPHABET ▾



Why ethics are important: consider the ‘privacy vs. security’ debate

Director of the FBI, James Comey, **called for “a regulatory or legislative fix” for technology companies’ expanding use of encryption to protect user privacy.** The post-Snowden pendulum has swung too far in one direction – in a direction of fear and mistrust. Justice may be denied because of a locked phone or an encrypted hard drive. Without a compromise homicide cases could be stalled, suspects could walk free, and child exploitation victims might not be identified or recovered”

“Privacy has never been an absolute right”



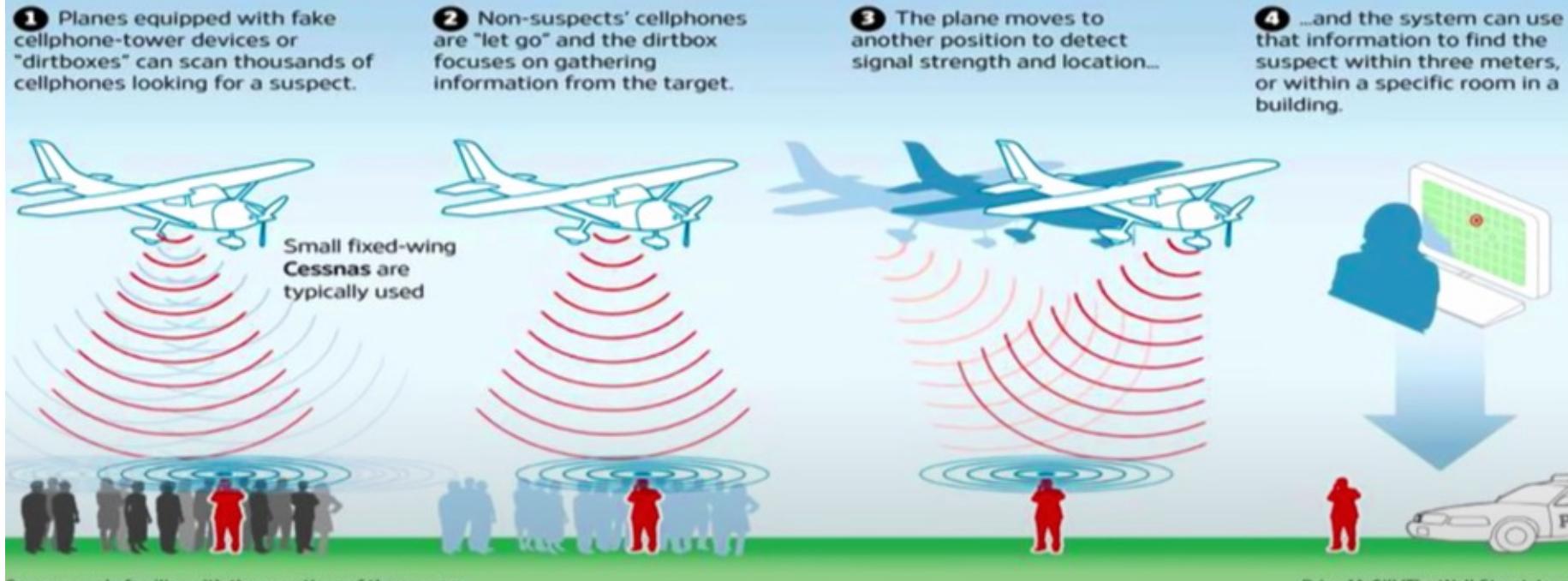
“I think it’s important to recognise that you can’t have 100 per cent security and also then have 100 per cent privacy and zero inconvenience”

US President Barack Obama



Almost every decision about technology usage now has ethical implications

Dirtboxes on a Plane | How the Justice Department spies from the sky



Source: people familiar with the operations of the program

Brian McGill/The Wall Street Journal



Brian McGill @brian_mcgill · 13h

The Justice Dept. is using 'dirtboxes' to scan cellphones in order to find criminal suspects. on.wsj.com/1xm0L9d

362

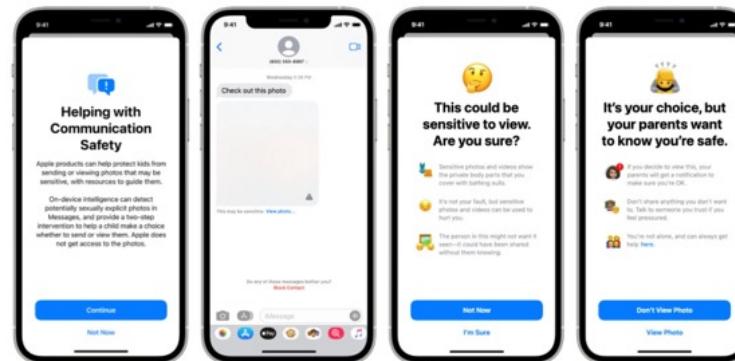
152

Twitter: @gleenhard





Apple delays the rollout of child-safety features over privacy concerns.



One Apple feature would allow parents to activate an alert when their children sent or received nude photographs in text messages. Apple

CSAM (Child Sexual Abuse Material)

<https://www.nytimes.com/2021/09/03/business/apple-child-safety.html>



Technology is progressing at exponential 'warp speed' while our ethics, social contracts and laws remain linear



Technology has seemingly limitless potential to improve our lives
– but should humans themselves **become technology?**



Technology has seemingly limitless potential to improve our lives
– but should humans themselves **become technology?**





Jibo - 'the world's first family robot'

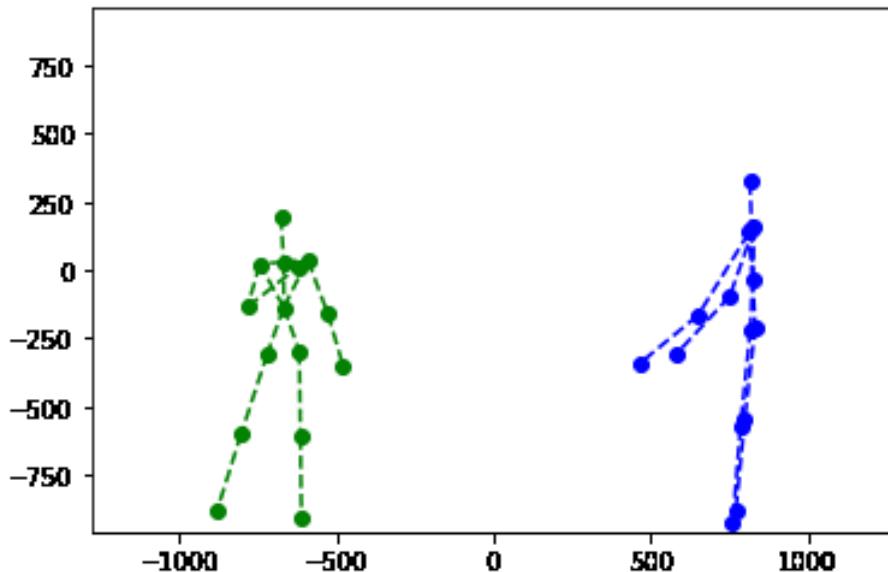
She has been working for two years on a product called Jibo, a machine she describes as a family robot because she envisions a day when every household will have one on the kitchen counter, playing an active and cheerful role in not only managing daily life but in making it better. Popular Mechanics

Social robots are not a tool, she argues. They are a partner. Tools force you to leave the moment. Jibo, she argues, will allow you to access all that information and technology while you stay in the moment—while you stay in

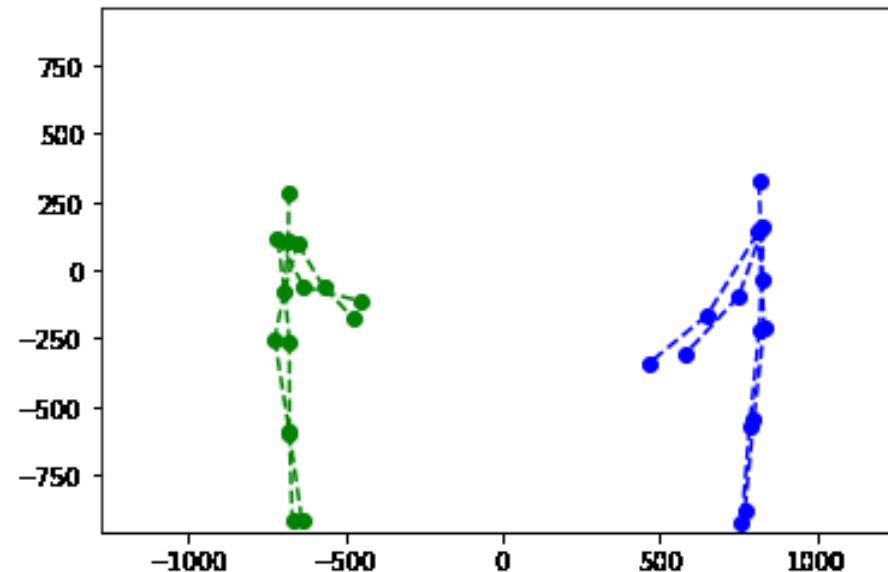
The image shows a woman in a kitchen, wearing an apron, standing at a counter. A white, spherical social robot, the Jibo, sits on the counter next to her. The background shows a window with flowers and a tiled wall.



Handshaking (normal)



Punching (adversarial)



One of my research:
Adversarial interaction attack: fooling AI to
misinterpret human intentions



1. Humans should not become technology
2. Humans should not be subject to dominant control by AI / AGI entities
3. Humans should not fabricate new creatures by augmenting humans or animals



Ethics in IT workplace



Ethical Problem Solving



Image from: <https://www.wikihow.com/Solve-Ethical-Issues>



What is an ethical dilemma?



What is an ethical dilemma?

Three conditions must be present for a situation to be considered an ethical dilemma:

1. An individual, the “agent”, must make a decision about which course of action is best. Situations that are uncomfortable but don’t require a choice, are not ethical dilemmas;
2. There must be different courses of action to choose from;
3. No matter what course of action is taken, some ethical principle is compromised, i.e. **there is no perfect solution.**

Allen, K. (2012). What is an thical dilemma? *The New Social Worker*. available at http://www.socialworker.com/feature-articles/ethics-articles/What_Is_an_Ethical_Dilemma%3F/



What is an ethical dilemma?

The Trolley problem (电车难题)

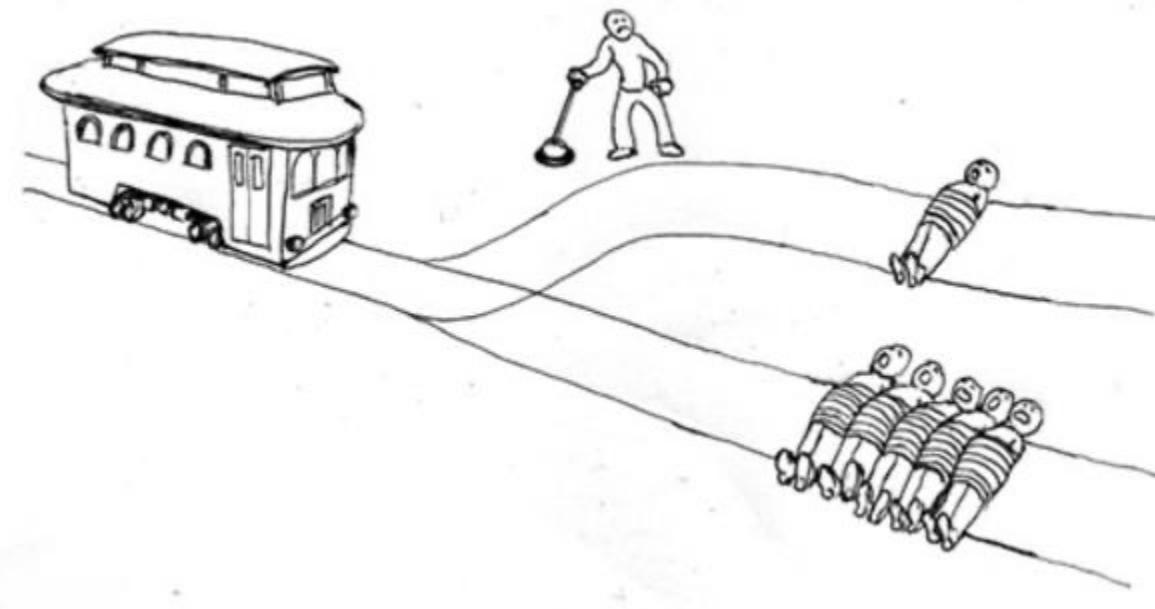


Image from: <http://knowyourmeme.com/memes/the-trolley-problem>



Case studies

- Data access
- Confidentiality
- Safety
- Trust
- Intellectual Property
- Privacy



Data access

- You are working for a Financial services industry company doing auditing / financial transaction management
- You get a database of a company' s finances showing its in big trouble and about to go bust, losing all its shareholders money
- You realize your elderly parents have invested all their life savings in this company and will go broke when this happens
- What do you do???



Confidentiality

- You work for a medical diagnostic device company
- You see information about real patients and their tests
- You recognize the name of your sibling's partner
- They have a very serious communicable disease you think your sibling hasn't been told about
- What would you do???



Safety

- You are working for a company developing self-driving cars
- You have to decide what to do in the critical scenario of the car either (i) hit another car or (ii) hit a pedestrian
- What do you program the car's AI to do?
- What could/should happen in such a scenario?
- What do human drivers do now?
- Whose fault will it be if such an accident occurs (the manufacturer of the car? The programmers? The person in the car? The other car? The pedestrian e.g. if walking on the freeway?)



Billing your IT work – the truth or not?

- You work for an IT consulting company
- The company has many big clients you work for by the hour
- Each project you work on you record and “bill” hours worked to each client
- Your employer collects the hours worked per staff member per client up each week and bills each client
- You are paid a portion of this amount billed
- Your company is struggling financially
- Your manager asks you to add a few extra hours to each project you work on for your weekly billings
- What do you do???



Intellectual Property

- You are working for a company developing new mobile phone apps
- The company has a number of clients
- You and two friends decide that you could do a much better job of the applications working for yourselves in a new company
- You copy the most interesting bits of the designs and code to an off-site server
- You copy the customer database
- You start your new company, develop new apps, and approach the clients to sell it
- What do you think will happen? Why?
- Your old boss sues you and your new company for extensive damages and restraint of trade



Privacy

- You work for a company building computer games using technologies such as Kinect, phones, Virtual Reality and Augmented Reality.
- The company captures information about gamers e.g. demographics, playing times, what games they are playing, background conversations, etc.
- Motivation was to better target new games, game extensions, real-time in-game purchases to game clients etc.
- Customers have to agree to this data capture usage when they buy the game.
- Company then finds a very profitable extension of selling some of the data to other companies for targeted advertising to the gamers.
- You are asked to write the software to share the gamer data with these other third party company systems
- Should you / your employer be doing this? Who will get prosecuted for breach of privacy and data laws??



Mobile devices – monitoring, data usage consent

- You manage an IT team that spends a lot of its time visiting client offices around the city/state.
- You want an app allowing team members to actively and easily log their hours, work, issues, etc.
- You are concerned about team members safety and security. Therefore, you also want “passively” monitor their whereabouts using GPS.
- On obtaining the data after deployment of the app you find one team member spends a lot of “work” time in the casino, while logging work activities there. The casino operators are not a client of your company.
- But – you didn’t get staff agreement to allow you to use the data in this way from the app.
- How do you handle this situation as the IT team manager?
- When you become a manager of other staff, you will get lots and lots of challenging Human Resources/ ethical issues to deal with ☺



And lots and lots more

...



Social Responsibility



Image from: <http://gateleyplc.com/corporate-social-responsibility/>



On Social Responsibility in IT

- As IT Professionals we have a lot of responsibilities to the community, stakeholders, and each other.
- For example: safety and security of systems, maintaining privacy and confidentiality, protection of critical infrastructure, intellectual property, plagiarism, ethical behavior...
- What would you do? Why? How can we know the “right” thing to do??



Example – Ariane 5 Rocket

1996年阿丽亚娜5火箭首飞失败，在发射40秒后**爆炸**



Ariane 5**爆炸**是历
史上最昂贵的软件
错误之一

https://www.youtube.com/watch?v=gp_D8r-2hwk



Example – Ariane 5 Rocket

- 4 June 1996
- ~40 seconds into launch
- altitude of ~3700m
- launcher veered off path and broke up
- then exploded
- ~\$500 million uninsured (maiden flight)
- Un-manned



What happened

- Technically:
 - Data conversion from a 64-bit floating point was too large for the target 16-bit signed integer value
 - Data conversion was not protected
- Causes:
 - The software module in question actually served no useful purpose after launch!
 - Was a carry over from Ariane 4
 - Operand error occurred because Ariane 5 built up a horizontal velocity much more quickly than Ariane 4
- **Whose fault was this? What could/should have been done?**



2016 Australian Census debacle

- System too slow / didn't respond
 - Had to do it – people panicked
 - Administrators thought under cyber-attack – shut it down
 - Provider didn't do sufficient scalability testing
 - Political, economic, social fall-out
-
- **Who is responsible? What could/should have been done?**



C U Next Week!

Course page:

<https://trustworthymachinelearning.github.io/>

Textbook:

下载链接:

Email: xingjunma@fudan.edu.cn

Personal page: www.xingjunma.com

Office: 江湾校区交叉二号楼D5025

