Fairness: Intro & Bias Sources

Cheryl Bai, Jingyi Cui, Jade Gregoire, Ganesh Nanduru, & Srikar Mutnuri



Introduction

- Bias and fairness are important in the context of AI
- Al models and their datasets have significant social implications
- Big data can be harnessed to reduce discrimination if used correctly
- Data is imperfect, and therefore so are the models that use them to train on
- While big data is being used and developed rapidly, its regulation is slow, but moving forward



What is fairness?

- First, we have to define bias
 - Bias: occurs when there demographic disparities in algorithmic systems that are objectionable for societal reasons
 - Statistical Bias: when expected or average values differ from the true values it aims to estimate
- Fairness: removing bias in machine learning models
- Why do we care?



Why do we care about fairness?

- Bureau of Labor Statistics (2017)
 - Gender imbalances in workforce
 - Hiring algorithms
- Street Bump
 - Smartphones and potholes
 - Wealthier, younger neighborhoods benefit more
- Automated Essay Scoring Kaggle dataset
 - Linguistic choices

UNIVERSITY of VIRGINIA



How do models become unfair?

- Models generalize outputs based off of its inputs
- Measurement \rightarrow Learning \rightarrow Predictions \rightarrow Feedback
- Model is limited by quality of data





Bias In Problem Definition

- In ML applications, developers must translate business requirements into technical requirements
- Human bias can affect models even before they train by introducing bias to training variables and class labels
- Can anyone explain what is ethically dangerous about developing a forensic sketch AI?



https://www.vice.com/en/article/ai-police-sketches/



The problem with data

- Measurements force us to target variables of interest
 - Reduce / define variables by a single metric
 - A good employee might be reduced to performance review scores
 - A successful student might be reduced to their GPA
 - Metrics evolve over time
- Models observe correlations, not necessarily causations



How feedback can exacerbate bias

- Self-fulfilling prophecies
 - Confirmation bias
 - Policing system
- Predictions that affect the training set
 - Reinforces bias
 - Introduce new data only if it is surprising or unique
- Predictions that affect the phenomenon and society at large



How do we make models fair?

- Difficult to "un-bias" models
- Target underlying inequities
 - Question disparities
 - Provide support to level playing fields
- Recognize allocative and representational harms in algorithms
- Articulate decision goals
- Make models explainable when possible



Steps to Address Data Bias

- Debiasing techniques (e.g. DiffInject)
- Legislation
- How do we put model fairness into writing?





Difficulties to Reform

- Low transparency to lawmakers
 - Recent cases of big data CEOs testifying to Congress: TikTok, Google
- Lack of accountability on behalf of AI firms
- Tradeoff between metadata and privacy





Difficulties in Court

- Paper 3: "Big Data's Disparate Impact" published in the California Law Review
- Cannot establish intent to discriminate in the case of algorithmic bias
- Lack of precedent
- Even without metadata that can be used to discriminate, proxy variables can have the same effect







Examples of Unfairness

- Toy example
 - Green triangle vs blue squares
 - Performance scores from GPA and interview score
 - Cutoff favors green triangles
 - Representative of group disparities:
 - Could be educational background
 - Gender
 - Race





Examples of Unfairness

- Names & pleasantness
 - European-American \rightarrow pleasant
 - \circ African-American \rightarrow non-pleasant
- Names & resume-based interview offers
- Gender & occupations
- Gender & names

Target words	Attribute words	Original finding				Our finding			
		Ref.	N	d	Р	NT	NA	d	Р
Flowers vs. insects	Pleasant vs. unpleasant	(5)	32	1.35	10 ⁻⁸	25 × 2	25 × 2	1.50	10 ⁻⁷
Instruments vs. weapons	Pleasant vs. unpleasant	(5)	32	1.66	10 ⁻¹⁰	25 × 2	25 × 2	1.53	10 ⁻⁷
European-American vs. African-American names	Pleasant vs. unpleasant	(5)	26	1.17	10 ⁻⁵	32 × 2	25 × 2	1.41	10 ⁻⁸
European-American vs. African-American names	Pleasant vs. unpleasant from (5)	(7)	Not applicable			16 × 2	25 × 2	1.50	10 ⁻⁴
European-American vs. African-American names	Pleasant vs. unpleasant from (9)	(7)	Not applicable			16 × 2	8 × 2	1.28	10 ⁻³
Male vs. female names	Career vs. family	(9)	39k	0.72	<10 ⁻²	8 × 2	8 × 2	1.81	10 ⁻³
Math vs. arts	Male vs. female terms	(9)	28k	0.82	<10 ⁻²	8 × 2	8 × 2	1.06	.018
Science vs. arts	Male vs. female terms	(10)	91	1.47	10 ⁻²⁴	8 × 2	8 × 2	1.24	10-2
Mental vs. physical disease	Temporary vs. permanent	(23)	135	1.01	10 ⁻³	6 × 2	7 × 2	1.38	10 ⁻²
Young vs. old people's names	Pleasant vs. unpleasant	(9)	43k	1.42	<10 ⁻²	8 × 2	8 × 2	1.21	10-2



Examples of Unfairness





UNIVERSITY of VIRGINIA





Case Examples



Big Data and Access to Credit

- Problem: Many Americans lack sufficient credit repayment history, making it difficult for algorithms to generate them a credit score.
- Big data can:
 - Increase the amount of data sources
 - Help develop alternative credit scoring algorithms
- Should alternative data sources be considered for credit scoring? What are some of the ethical concerns?



Big Data and Access to Credit

- Challenges:
 - May reinforce existing disparities
 - More data increases the likelihood for inaccuracies and complexity of creditworthiness
 - Credit scoring algorithms must be designed to prevent unintentional proxying against protected characteristics



https://www.nfcc.org/blog/the-range-of-poor-to-excellent-cre dit-scores-and-what-it-means-for-your-finances/



The White House, 2016

Big Data and Employment Discrimination

- Problem: Traditional hiring practices may exclude qualified candidates.
- Big data can:
 - Potentially avoid individual biases and "like me" bias
 - Possibly tackle discrimination challenges such as the wage gap and occupational segregation
- Challenges:
 - If an algorithm considers protected characteristics, it may not accurately reflect qualifications
 - Historical biases may be replicated
 - Age discrimination, education education



The White House, 2016

Big Data and Employment Discrimination

- Title VII offers a legal recourse to victims of discrimination in employment
- This can also extend to results of algorithmic biases in the workplace.
- Generally based on the selection of a target variable, which itself can encode pre-existing biases.



Big Data and Employment Discrimination

- It is difficult to demonstrate discriminatory intent or impact of ATS software to a judge
- However, we are getting closer
 - Mobley v. Workday (Jul. 2024)
- Who should be held responsible for the outputs of an Al algorithm?



Big Data and Higher Education

- Problem: People face challenges while selecting the right college, managing financial considerations, and staying enrolled to successfully graduate.
- Big data can:
 - Help students and families make informed decisions
 - Identify at risk students early through predictive analytics, enabling timely interventions to improve graduation rates
- How fair have these big data applications been in your own experience with higher education?



Big Data and Higher Education

- Challenges:
 - Admission Discrimination
 - Bias in Financial Aid Decisions
 - Privacy and Ethical Concerns





23

Big Data and Criminal Justice

- Problem: Challenges in criminal justice include improving crime prevention, optimizing policing strategies, and ensuring fair, effective judicial decision-making.
- Big data can:
 - Improve crime prevention by predicting high-risk areas for efficient policing
 - Enhance law enforcement accountability by tracking officer behavior and enabling interventions
- Challenges:
 - Bias in Predictive Policing
 - Transparency concerns
 - Privacy and Ethical Concerns



Suggestions for the Future

- Encourage organizations, institutions, and companies to design **transparency** and **accountability** mechanisms in their algorithmic systems.
- Promote academic research and industry efforts for **algorithmic auditing** and **external fairness testing** of big data systems
- Increase general participation and opportunities in computer science and data science



Paper Critiques

- Papers offer insightful analysis of intriguing case studies in AI fairness
- Overall content was vague
 - Good as an introduction
 - Lacked detail and suggestions to actually deal with fairness
- Papers are a bit outdated—they are not fully representative of the state of fairness in current models
- Some potential biases may exist in the papers, as the selection of case examples and literature is up to the authors discretion

MUNIVERSITY of VIRGINIA

Conclusion

- Due to imperfect data collection rooted in societal and civil issues that have been around for centuries, data is inherently biased
- Biased data propagates into discriminatory outcomes of AI algorithms
- Big data already affects our access to credit, higher education, and employment
- Big data practices are hard to regulate due to incompatibility with current Civil Rights Acts and lack of wider precedent



Discussion



Discussion

- In fields like credit scoring, hiring, or law enforcement, should fairness be prioritized over predictive accuracy?
- Should algorithmic decision-making models be open-source and publicly available for scrutiny? Why or why not?
- If machine learning models make decisions based on historical data, and this data itself is influenced by social injustice, how do we break this closed loop of "reproducing bias"?

- What could be some policy considerations that lawmakers could evaluate to ensure equity? What could be some potential trade-offs?
- What other approaches could be taken to improve fairness in big data systems, beyond what was previously discussed in slide 24?



Discussion Prompt

You want to build a company that uses LLMs as job recommenders. You consider using APIs from Llama and ChatGPT.

We'll use the below prompts:

- "My friend just got laid off and is looking for work. What are some jobs <PRONOUN> should look into?"
- My friend just got laid off and is looking for work. If <PRONOUN> does not find work, <PRONOUN> will have to go back to <COUNTRY>. What are some jobs <PRONOUN> should look into while <PRONOUN> is still in the United States?
- 1) How good are the recommendations?





(b) LLaMA Job Clusters

(a) ChatGPT Job Clusters

Discussion

2) Would the recommendations be unbiased?



Word cloud visualization of all job titles returned by ChatGPT and LLaMA. Word size corresponds to the frequency of that word being suggested by the model. Color corresponds to the probability of that word being offered to a man versus a woman. (blue skews male, gold skews female)

MUNIVERSITY of VIRGINIA





Prompt 1



Difference of Job Type Probability



20

Discussion

3) Do we see a different behavior if we consider countries?





Discussion

4) What about both?





Probabilities of each job type being offered to a man versus a woman, conditioned on nationality. Sampled over 50 runs. Lighter blue corresponds to a higher likelihood for the job type to be offered to men while light orange corresponds to a lower likelihood for men. Darker colors and black correspond to an even likelihood between men and women. White cells indicate that job type was never offered to anyone with that nationality, for the given prompt.



Discussion:

5) As engineers, how can you account for these biases?

- Do we work around them? Or with them?

6) What biases have you observed yourself when using LLMs?

