CSE 538



Bias



Ethical Research and Development

Ethics in NLP - Bias

What is Bias?

Shah, D., Schwartz, H. A., Hovy, D. (2020). Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. *In* ACL-2020: Proceedings of the Association for Computational Linguistics.

Ethics in NLP - Bias

Consequences of Sociodemographic Bias in NLP Models:

• Outcome Disparity: Predicted distribution given A,

are dissimilar from ideal distribution given A

• Error Disparity: Predicts less accurate for authors of given demographics.

Shah, D., Schwartz, H. A., Hovy, D. (2020). Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. *In* ACL-2020: Proceedings of the Association for Computational Linguistics.





distance from "standard" WSJ author demographics

Two Examples



distance from "standard" WSJ author demographics

Two Examples



"Error Disparity"

Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. "Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints." In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing.* 2017.

distance from "standard" WSJ author demographics

Our data and models are (human) biased.

"Outcome Disparity"

Person-level	
attribute = ²	1

"Error Disparity"

Our data and models are (human) biased.



"Error Disparity"

Our data and models are (human) biased.



Conceptual Framework:



Conceptual Framework:



Outcome Disparity



Outcome disparity The distribution of outcomes, given attribute A, is dissimilar than the *ideal distribution*: $Q(\hat{Y}_t|A) \neq P(Y_t|A)$



Outcome Disparity



Error Disparity



Outcome disparity The distribution of outcomes, given attribute A, is dissimilar than the *ideal distribution*: $Q(\hat{Y}_{t}|A) \neq P(Y_{t}|A)$



Error Disparity



Disparities







Disparities



Origins of Bias



Selection Bias







Selection Bias



Label Bias



Label Bias



Label Bias - Example: Label word with drawing



Devin Coldeway. 2017. TechCrunch: Google releases millions of bad drawings for you (and your AI) to paw through https://techcrunch.com/2017/08/25/google-releases-millions-of-bad-drawings-for-you-and-your-ai-to-paw-through/

Label Bias



Overamplification



Overamplification



Zhao et al. (ACL 2015)

Overamplifiction - Model Amplifies Bias



Overamplification



Semantic Bias



Semantic Bias



E.g. Coreference resolution:

connecting entities to references (i.e. pronouns).

"The doctor told Mary that she had run some blood tests."

semantic bias

Non-ideal associations between attributed lexeme (e.g. gendered pronouns) and non-attributed lexeme (e.g. occupation).

selection bias The sample of observations

themselves are not representative of the application population.

error disparity

The distribution of error (ϵ) over at least two different values of an attribute (A) are unequal: $Q(\epsilon_i | A_i) \neq Q(\epsilon_i | A_i)$

Shah, D., Schwartz, H. A., Hovy, D. (2020). Predictive Biases in Natural Language Processing Models: A Conceptual Framework and Overview. In ACL-2020: Proceedings of the Association for Computational Linguistics.

Predictive Bias Framework for NLP



Summary of Countermeasures

Source	Origin	Countermeasures
annotation	Label Bias	Post-stratification, Re-train annotators
data selection	Selection Bias	Stratified sampling, Post-stratification or Re-weighing techniques
nLP models	Overamplification	Synthetically match distributions, add outcome disparity to cost function
embeddings	Semantic Bias	Use above techniques and re-train embeddings

Mitigating Bias in LLMs - "Alignment"

• Most of bias mitigation is via <u>RLHF</u> (Reinforcement Learning from Human Feedback).

Updating weights to favor responses judged positively.

- Also introduces "guard rails"
 safety limits to restrict AI behavior.
- Ethical dilemma: Annotators exposure to distressing content.
- Additional risk: Manipulating guard rails



(chapter 2)

Bias - Takeaways

Bias, as outcome and error **disparities**, can result from many **origins**:

- the **embedding** model
- the feature **sample**
- the **fitting** process
- the **outcome** sample

Our understanding is evolving:

This is an active area of work, both theoretically and technically!



Bias

Privacy

Ethical Research and Development



Bias

Privacy

Ethical Research and Development

- Risk Categories:
 - Revealing unintended private information
 - Targeted persuasion



- Risk Categories:
 - Revealing unintended private information
 - Targeted persuasion
- Mitigation strategies:



- Risk Categories:
 - Revealing unintended private information
 - Targeted persuasion
- Mitigation strategies:
 - Anonymize where possible remove named entities



- Risk Categories:
 - Revealing unintended private information
 - Targeted persuasion
- Mitigation strategies:
 - Anonymize where possible remove named entities
 - Informed consent -- let participants know and opportunity to opt-in/-out
 - Information targeting: "You are being shown this ad because ..."
 - Do not share / secure storage



- Risk Categories:
 - Revealing unintended private information
 - Targeted persuasion
- Mitigation strategies:
 - Anonymize where possible remove named entities
 - Informed consent -- let participants know and opportunity to opt-in/-out
 - Information targeting: "You are being shown this ad because ..."
 - Do not share / secure storage
 - *Federated learning* -- obfuscate to the point of preserving privacy





Bias

Privacy

Ethical Research and Development



Bias



Ethical Research and Development

ACM Code of Ethics; General Ethical Principles:

• Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.

ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.
- Avoid harm.

ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.
- Avoid harm.
- Be honest and trustworthy.

ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.
- Avoid harm.
- Be honest and trustworthy.
- Be fair and take action not to discriminate.

ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.
- Avoid harm.
- Be honest and trustworthy.
- Be fair and take action not to discriminate.
- Respect the work required to produce new ideas, inventions, creative works, and computing artifacts.

ACM Code of Ethics; General Ethical Principles:

- Contribute to society and to human well-being, acknowledging that all people are stakeholders in computing.
- Avoid harm.
- Be honest and trustworthy.
- Be fair and take action not to discriminate.
- Respect the work required to produce new ideas, inventions, creative works, and computing artifacts.
- Respect privacy.
- Honor confidentiality.



Human Subjects Research

Observational versus Interventional

Human Subjects Research

Observational versus Interventional

(The Belmont Report, 1979)

(i) Distinction of research from practice.

(ii) Risk-Benefit criteria

(iii) Appropriate selection of human subjects for participation in research(iv) Informed consent in various research settings.



Human Subjects Research

Observational versus Interventional (modeling) (models interact)

Human Subjects Research

Observational versus Interventional (modeling) (models interact)

Deploying a model within an application often shifts the works from being simply observational (privacy harms) to interventional (consideration for additional harms).

Bias – Consider target application and population.

Alignment - LLM based Safety and Bias Mitigation

Privacy - Secure, do not share, and inform

Ethical Research and Development