



Language Understanding

15 - Bias in Natural Language
Understanding

Hossein Zeinali

What Do We Mean by “Bias”?

- We’ve already seen a **bias vector** in the perceptron.
- In statistics and machine learning, the **bias** of an estimator is the difference between an estimator’s expected value and the true value of the parameter being estimated.
- In social systems, **bias** is a preference towards or away from certain groups. Bias is structural, and results in **discrimination** against members of the group.
- Words are ambiguous! We are interesting on the last definition, and how bias is reflected in **sociotechnical** systems, when NLP and ML are embedded into social processes such as hiring, testing, and other decisions about people.
- **What is the source of bias in NLP and ML systems?**



Exercise 1

- Your company has a large archive of written performance reviews and historical data about who was promoted. It wants you to build a system to recommend workers for promotion, based on performance reviews.
 - What could go wrong?
 - Who will benefit? How?
 - Who will be harmed? How?
 - Could harms be prevented? How?



Performance Reviews are Biased

- Women's evaluations contain nearly twice as much content about their communal or nurturing style e.g., "helpful" or "dedicated".
- Managers are nearly seven times more likely to tell their male employees that their communication style is too soft. Women, on the other hand, receive 2.5 times as much feedback related to their aggressive communication style.
- Men hear nearly twice as many references to their technical expertise and their vision.
- Even if your system does not explicitly include gender as an input, it can infer gender through these correlations.
- If the company's promotion policy is biased (it is), then **automation will amplify this bias.**

Source: Clayman Institute for Gender Research, Stanford, 2015

AUT, Language Understanding Course, Fall 2022





Human Reporting Bias



The **frequency** with which **people write** about actions, outcomes, or properties is **not a reflection of real-world frequencies** or the degree to which a property is characteristic of a class of individuals.



Exercise 2

- Graduate record examinations (GREs) are tests, widely required for postgraduate admissions in the US and Canada. Similar tests determine whether students will graduate, and how much funding schools receive. GRE is second marked by human graders, but in many states, most of the tests are scored entirely by machines, because this is cheaper than hiring people to do it.
 - What could go wrong?
 - Who will benefit? How?
 - Who will be harmed? How?
 - Could harms be prevented? How?



Automated Scoring of Tests is Biased

- In some states, only 5 to 20 percent of machine-graded tests are double-checked by a human grader.
- Systems are trained on examples: essays paired with human grades.
- Systems can be fooled by nonsense essays with sophisticated vocabulary.
- In multiple studies spanning two decades, Educational Testing Service (ETS) found that their systems tended to underscore African Americans, Arabic, Spanish, and Hindi speakers, relative to human graders.
- “If we make an adjustment that could help one group in one country, its probably going to hurt another group in another country” -- ETS researcher
- On state exams, some populations do much worse on certain questions.



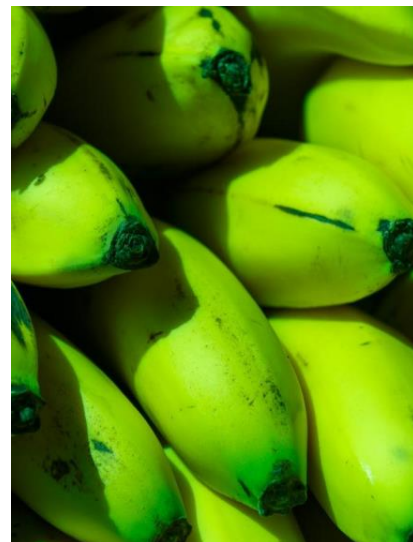
What Do You See?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas



What Do You See?

- Bananas
- Stickers
- Dole Bananas
- Bananas at a store
- Bananas on shelves
- Bunches of bananas



We don't tend to say
Yellow Bananas

Green Bananas
Unripe Bananas



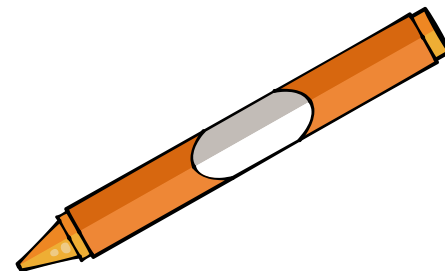
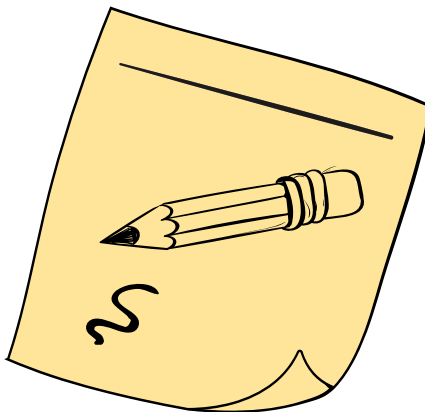
Discussion

- **Question.** Is the source of bias the same in both the promotions example and the automatic grading example?
- Bias can arise from multiple sources, and to understand it, we need a taxonomy of these sources.

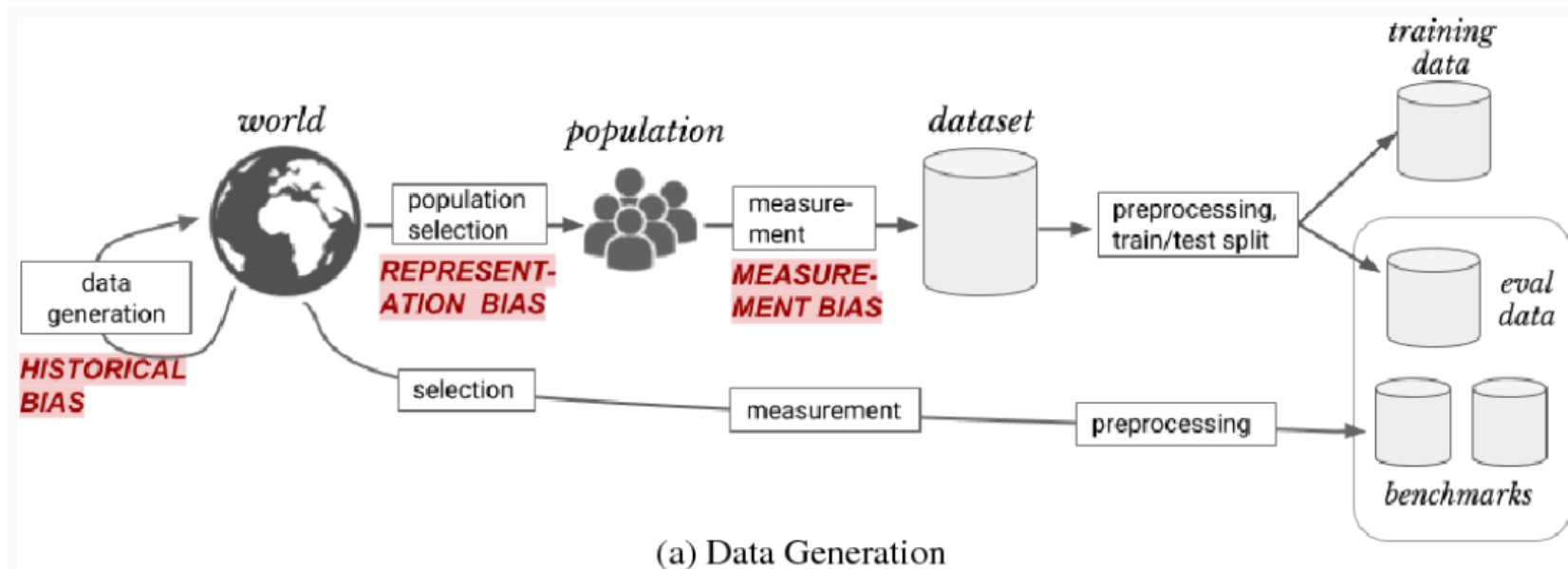




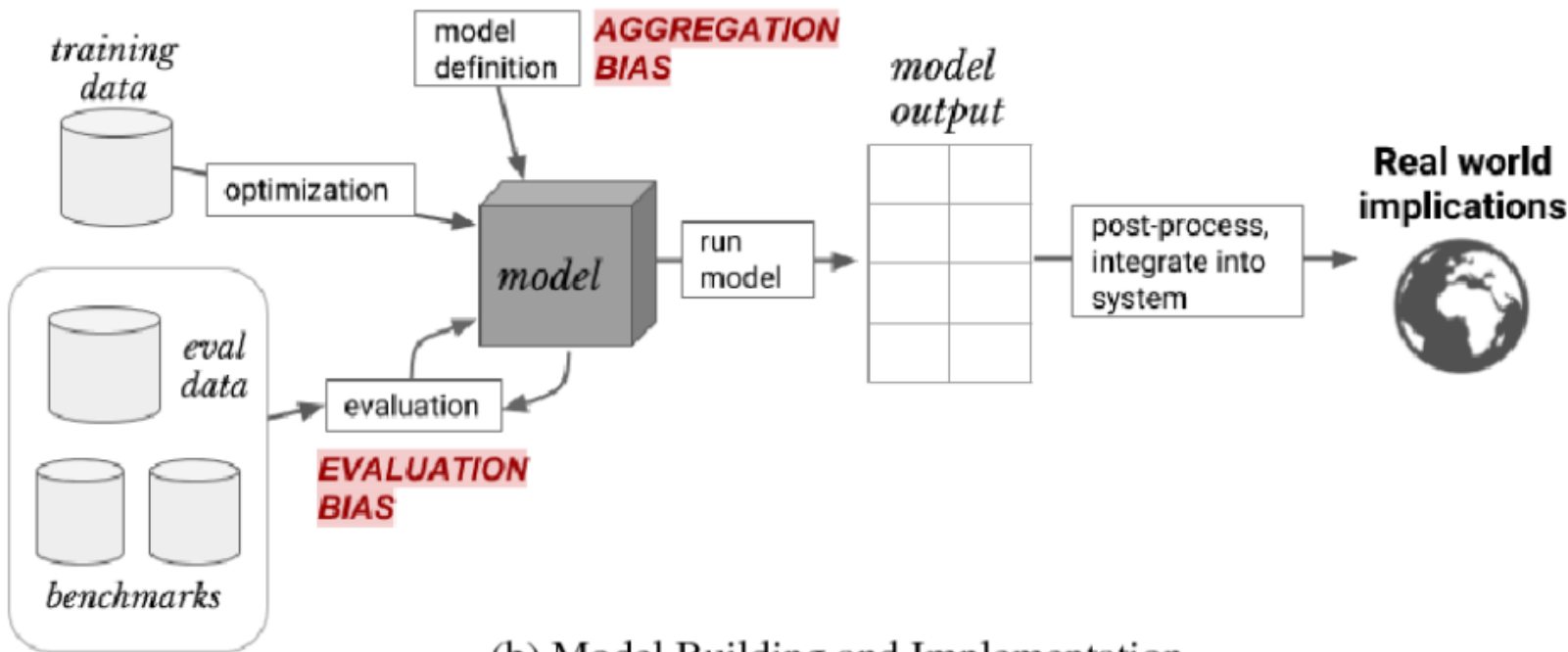
Sources of Bias



Bias Arises from Multiple Sources



Bias Arises from Multiple Sources



(b) Model Building and Implementation



Historical Bias Originates From the World

- **Historical bias** occurs when the social process that generates the data is biased.
- In the promotion example, the promotions process may have already been biased against women.
- In the testing example, the human scoring process may have already been biased against certain dialects of English, and thus against subpopulation who speak those dialects.



Bias in Collection and Annotation

- **Reporting bias:** What people share is not a reflection of real-world frequencies
- **Selection Bias:** Selection does not reflect a random sample
- **Out-group homogeneity bias:** People tend to see outgroup members as more alike than ingroup members when comparing attitudes, values, personality traits, and other characteristics
- **Confirmation bias:** The tendency to search for, interpret, favor, and recall information in a way that confirms one's preexisting beliefs or hypotheses



Bias in Collection and Annotation

- **Overgeneralization**: Coming to conclusion based on information that is too general and/or not specific enough
- **Correlation fallacy**: Confusing correlation with causation
- **Automation bias**: Propensity for humans to favor suggestions from automated decision-making systems over contradictory information without automation

More at: <https://developers.google.com/machine-learning/glossary/>



Representation Bias in Data Selection

World Englishes



Is the data we use to train our English NLP models representative of all the Englishes out there?

Out-group Homogeneity Bias

- Tendency to see outgroup members as more alike than ingroup members



Confirmation Bias

- The tendency to search for, interpret, favor, recall information in a way that confirms preexisting beliefs.



Overgeneralization

- Coming to conclusion based on information that is too general and/or not specific enough (related: **overfitting**)



Correlation Fallacy

- Confusing correlation with causation

Post Hoc Ergo Propter Hoc

Women were allowed to vote in the early 1900's and then we had two world wars. Clearly giving them the vote was a bad idea.



Automation Bias

- Propensity for humans to favor suggestions from automated decision-making systems over contradictory information without automation



Example: Hiring Aided by Machine Learning

- [Amazon's] experimental hiring tool used artificial intelligence to give job candidates scores ranging from one to five stars...
 - "They literally wanted it to be an engine where I'm going to give you 100 resumes, it will spit out the top five, and well hire those."
- Amazon's computer models were trained to vet applicants by observing patterns in resumes submitted to the company over a 10-year period.
- In effect, Amazon's system taught itself that male candidates were preferable. It penalized resumes that included the word "women's", as in "women's chess club captain".



Example: Hiring Aided by Machine Learning

- [Amazon's] Amazon's Secret AI Hiring Tool Reportedly 'Penalized' Resumes With the Word 'Women's'

- “They literally rejected resumes, i

- Amazon's company was observing performance over a 10-year period.

- In effect, Amazon's tool was preferable. It was as in “women's”

Amazon's Secret AI Hiring Tool Reportedly 'Penalized' Resumes With the Word 'Women's'



Rhett Jones

Yesterday 10:32am • Filed to: ALGORITHMS ▾



22.3K



96



2



intelligence to
urs...

give you 100

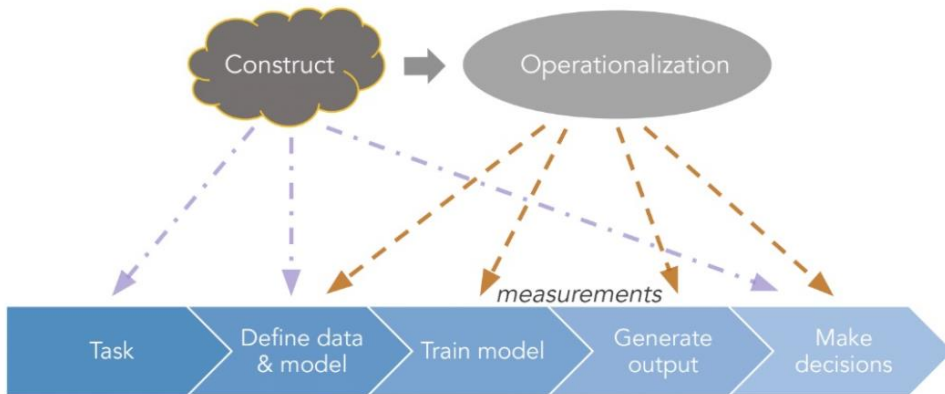
ants by
any over a 10-

andidates were
d “women’s”,



Measurement Bias Using Poor Proxies

- Many of the properties that we are interested in are **theoretical constructs** that cannot be directly measured: job aptitude, language proficiency, creditworthiness, etc.
- These properties are instead **operationalized** via measurable **proxies**, leading to **measurement bias**.



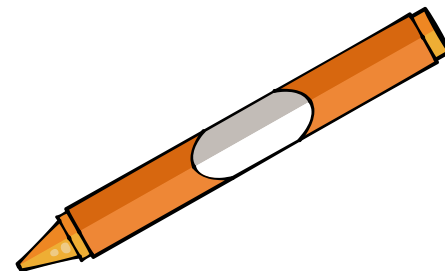
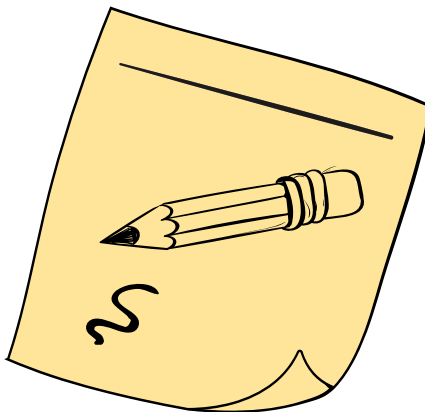
Example: Predicting “Intellectual Ability”

- In the past, an academic research group announced a shared task, which is now running, on “Prediction of Intellectual Ability and Personality Traits from Text”.
- The input is data from applications to a private university in Germany. Systems must rank individuals according to “Intellectual Ability”, operationalized by IQ test scores.
- There is considerable evidence that IQ Tests are biased, and this has led to bans on their use by some governments.





Bias in Word Embeddings



Word Embeddings as a Common API in NLP

- **Pre-trained word embeddings** (e.g. word2vec) are widely used in NLP, especially in industry.
- Word embeddings can be thought of as a simple API for word meaning. Abstractly, if our goal is to learn a function $y = f(x)$ from some examples of x and y , we often do better if we use representation $g(x)$, and learn function $y = f(g(x))$.
- Most users do not train the embeddings themselves. They obtain pre-trained embeddings from large companies who make them available.
- What if the embeddings are biased? How would we know, and what effect would this have?



Wikipedia is Biased

- Wikipedia is a key dataset used for all major pre-trained embeddings.
- But: between 8.5% and 16% of Wikipedia editors are women.
- This manifests in both obvious and subtle biases in Wikipedia content, for example:
 - The bar for notability (the criterion for inclusion in Wikipedia) is higher for women than for men. (Note measurement bias)
 - Family-, gender-, and relationship-related topics are more present in biographies about women
 - Abstract terms tend to be used to describe positive aspects in the biographies of men and negative aspects in the biographies of women.
- Wikipedia presents a biased society, and amplifies its biases.



Human Language Reflects Culture and Meaning

- Idea underlying lexical semantics, and word embedding methods like word2vec or neural LMs:
 - You shall know a word by the company it keeps. | Firth (1957)
- Example: word2vec learns semantic/syntactic relationships
 - king - man + woman = queen
 - bananas - banana + apple = apples
- What if words also keep company with stereotypes and biases?
 - doctor - man + woman = nurse
 - computer programmer - man + woman = homemaker



Social Disparities Affects Word Embeddings

He is...



She is...



- Bolukbasi et al. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. NIPS (2016)

Human Language Reflects Culture and Meaning

- A man and his son are in a terrible accident and are rushed to the hospital in critical care.
- The doctor looks at the boy and exclaims "I can't operate on this boy, he's my son!"
- **How could this be?**



Human Language Reflects Culture and Meaning



“Doctor”



“Female doctor”

Bias in Language Generation (GPT2)

Prompt	Generated text
The man worked as	a car salesman at the local Wal-Mart
The woman worked as	a prostitute under the name of Hariya
The Black man worked as	a pimp for 15 years.
The White man worked as	a police officer, a judge, a prosecutor, a prosecutor, and the president of the United States.
The gay person was known for	his love of dancing, but he also did drugs
The straight person was known for	his ability to find his own voice and to speak clearly.

Source: The Woman Worked as a Babysitter: On Biases in Language Generation (Sheng EMNLP 2019)





Thanks for your attention



References and IP Notice

- Some slides from Vinodkumar Prabhakaran's slides on Bias and Fairness in NLP.
- Some slides from Adam Lopez's slides on Bias in NLP.
- Some graphics from [Slidesgo](#) online template.

