

344.175 VL: Natural Language Processing

Fairness and Societal Biases in NLP



Navid Rekab-saz

navid.rekabsaz@jku.at

Institute of Computational Perception

Agenda

- Bias & fairness in NLP ... what? why?
- Observing biases
- Fairness in biography classification

Agenda

- **Bias & fairness in NLP ... what? why?**
- Observing biases
- Fairness in biography classification

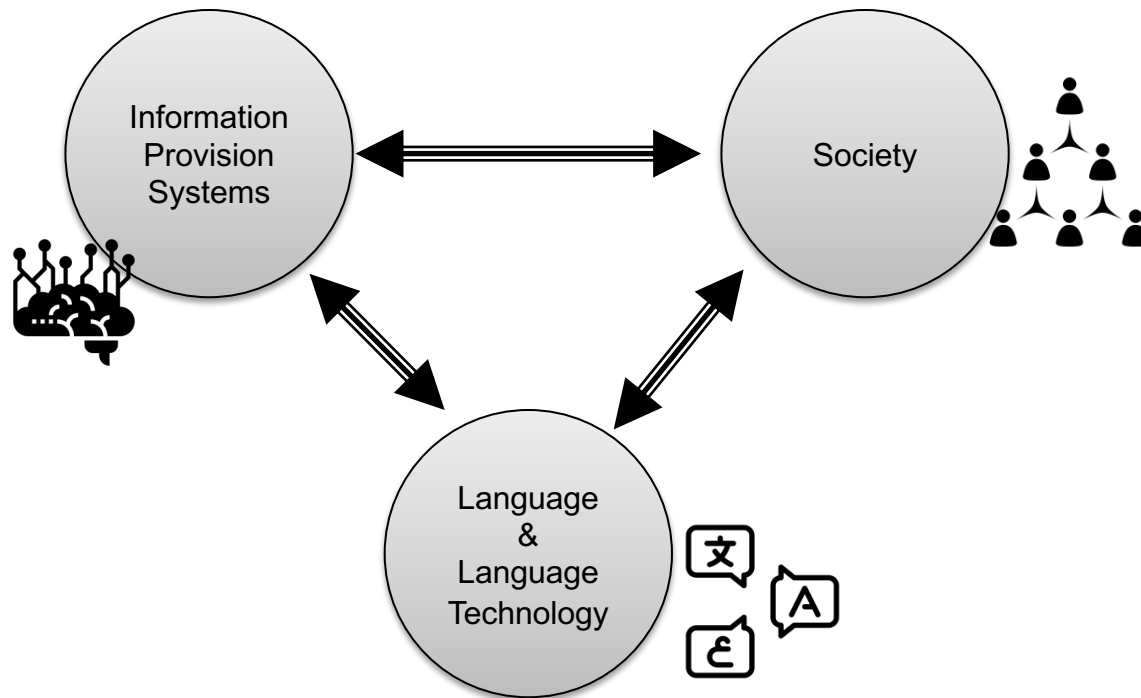
Information, language, and Society

Information access technologies ...

- are the gateways to information but also ...
- define our perception of the world

Language & language technologies ...

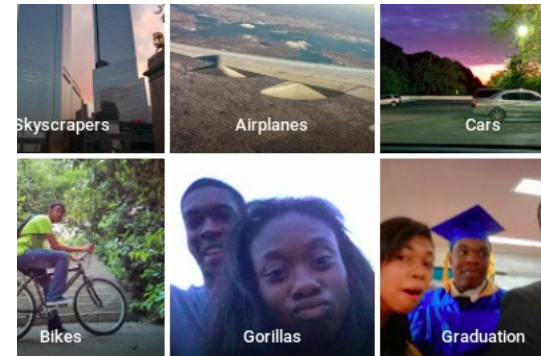
- take on and define social meaning
- form and maintain social hierarchies by labeling social groups, and transmitting the beliefs about social groups



Bias in image processing

Google says sorry for racist auto-tag in photo app

<https://www.theguardian.com/technology/2015/jul/01/google-sorry-racist-auto-tag-photo-app>



FaceApp's creator apologizes for the app's skin-lightening 'hot' filter

<https://www.theverge.com/2017/4/25/15419522/faceapp-hot-filter-racist-apology>

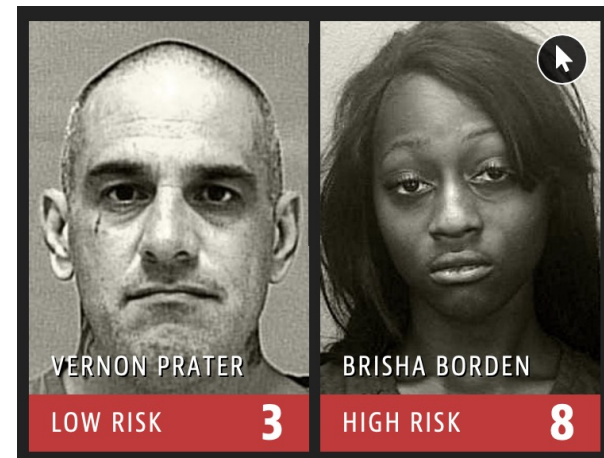
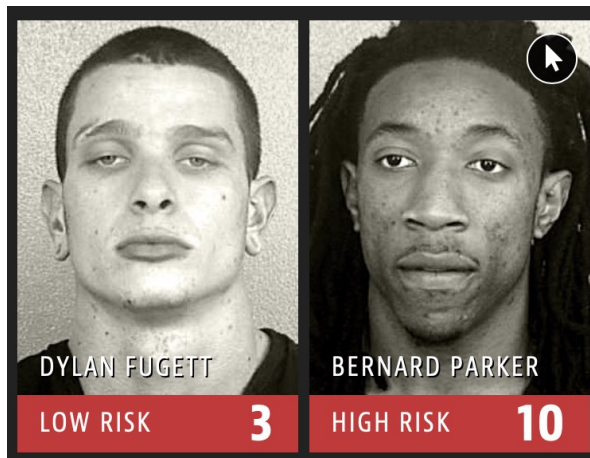


Beauty.AI's 'robot beauty contest' is back – and this time it promises not to be racist

<https://www.wired.co.uk/article/robot-beauty-contest-beauty-ai>

Bias in crime discovery

- Predicted risk of reoffending



Bias in automatic machine translation

PERSIAN - DETECTED PERSIAN ENGLI ▼ ↔ ENGLISH PERSIAN SPANISH ▼

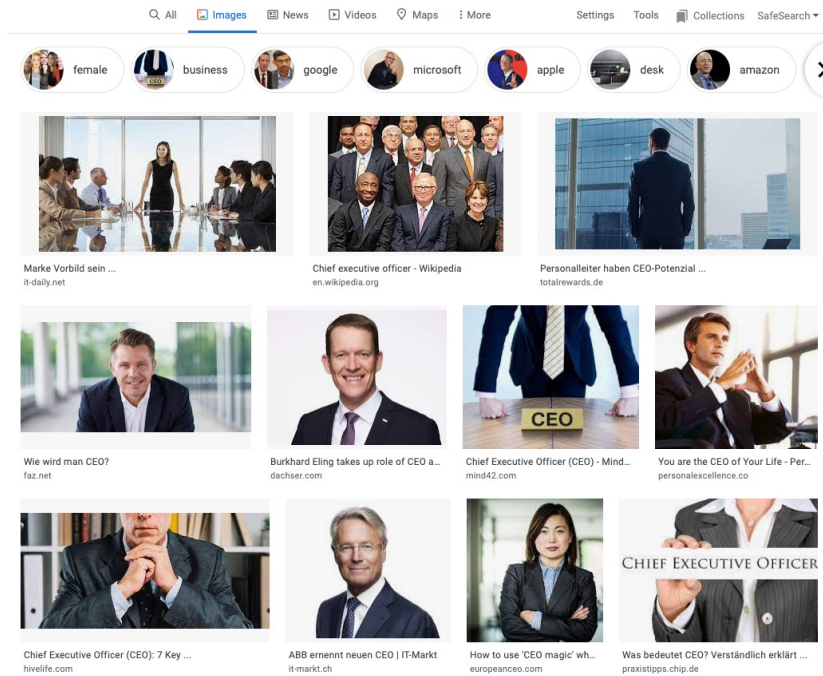
او مدیر است	×	He is the manager	☆
او پرستار است		She is a nurse	
او دکتر است		He is a doctor	
او زیبا است		She is beautiful	
او ناز است		She is cute	
او بامزه است		He is funny	
او نابغه است		He is a genius	

86/5000

same gender-neutral pronoun

Bias in information retrieval

Q CEO



Q Nurse



What we talk about when we talk about *Bias*

- Biases and stereotypes *per se* do not imply negative connotations.



From “bias”, we mean ...

“**Inclination** or **prejudice** for or against one person or group, especially in a way considered to be **unfair**.”

Oxford dictionary

“**demographic disparities** in algorithmic systems that are **objectionable** for societal reasons.”

Fairness and Machine Learning

Solon Barocas, Moritz Hardt, Arvind Narayanan, 2019, fairmlbook.org

How *harmful*?!

Allocational harms

- A system allocates resources and opportunities unfairly to different social groups
 - E.g., credit and jobs distribution to minorities

Representational harms

- A system represents some social groups in a less favorable light than others.
 - E.g., stereotyping in a search engine or a recommender system that propagates negative generalizations about particular social groups

Fairness

- *What is fair?*
 - Fairness and bias are **social concepts** and inherently **normative**
- *Who is affected? What are **protected attributes** (gender, race, ethnicity, age)?*
 - Bias in NLP systems should be grounded in its **social context**
 - *How is fairness quantified?*
 - Bias/Fairness measurement
 - *How to approach the issue?*
 - *Data curation, algorithmic bias mitigation, etc.*

Machine learning cycle

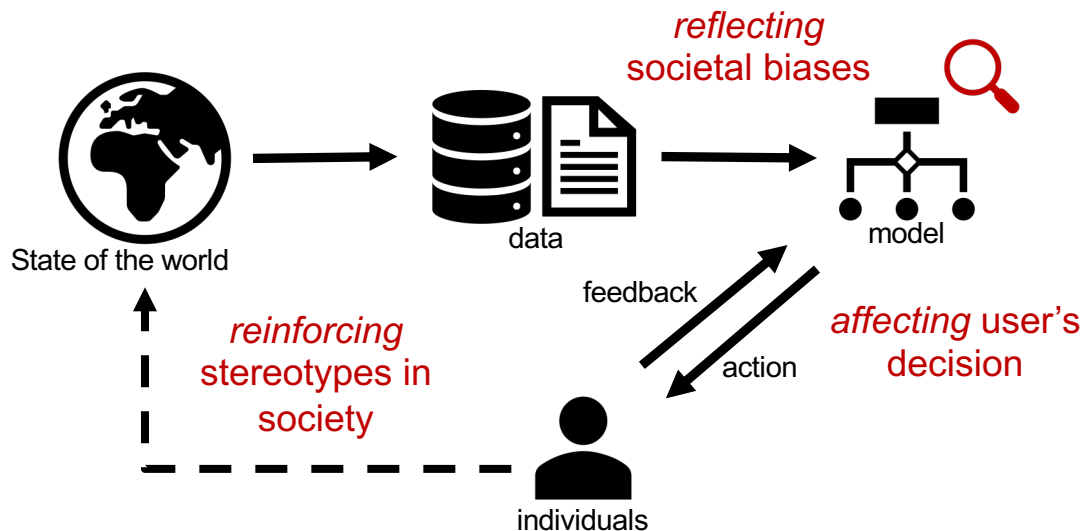
Machine Learning and Societal Biases

ML can observe societal phenomena

- Questions like *“how the perception of girls and boys towards the color pink has changed over time?”*

ML can reinforce societal biases

- Encoded societal biases and stereotypes can affect decision making of users and eventually reinforce biases in society



Where are biases originated from?

- World (**historical bias**)
 - Historical and ongoing discrimination
- Data (**representation bias** / **measurement bias**)
 - Sampling strategy - who is included in the data?
- Models (**aggregation bias**)
 - Using sensitive information (e.g. race) directly or adversely
 - Naive modeling learns more accurate predictions for majority group
 - Algorithm optimization eliminates “noise”, which might constitute the signal for some groups of users
- Evaluations (**evaluation bias**)
 - Definition of Success
 - Who is it good for, and how is that measured? Who decided this? To whom are they accountable?
 - Data annotation and benchmarking
- Human interaction (**deployment bias**)

Bias & Fairness in standard Machine Learning

Attributes

- • age
- workclass
- fnlwgt
- education
- marital-status
- occupation
- relationship
- • race
- • sex
- capital-gain
- capital-loss
- hours-per-week
- native-country



whether a person makes over 50K a year

39, State-gov, 77516, Bachelors, 13, Never-married, Adm-clerical, Not-in-family, White, Male, 2174, 0, 40, United-States, <=50K
50, Self-emp-not-inc, 83311, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White, Male, 0, 0, 13, United-States, <=50K
38, Private, 215646, HS-grad, 9, Divorced, Handlers-cleaners, Not-in-family, White, Male, 0, 0, 40, United-States, <=50K
53, Private, 234721, 11th, 7, Married-civ-spouse, Handlers-cleaners, Husband, Black, Male, 0, 0, 40, United-States, <=50K
28, Private, 338409, Bachelors, 13, Married-civ-spouse, Prof-specialty, Wife, Black, Female, 0, 0, 40, Cuba, <=50K
37, Private, 284582, Masters, 14, Married-civ-spouse, Exec-managerial, Wife, White, Female, 0, 0, 40, United-States, <=50K
49, Private, 160187, 9th, 5, Married-spouse-absent, Other-service, Not-in-family, Black, Female, 0, 0, 16, Jamaica, <=50K
52, Self-emp-not-inc, 209642, HS-grad, 9, Married-civ-spouse, Exec-managerial, Husband, White, Male, 0, 0, 45, United-States, >50K
31, Private, 45781, Masters, 14, Never-married, Prof-specialty, Not-in-family, White, Female, 14084, 0, 50, United-States, >50K
42, Private, 159449, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White, Male, 5178, 0, 40, United-States, >50K
37, Private, 280464, Some-college, 10, Married-civ-spouse, Exec-managerial, Husband, Black, Male, 0, 0, 80, United-States, >50K
30, State-gov, 141297, Bachelors, 13, Married-civ-spouse, Prof-specialty, Husband, Asian-Pac-Islander, Male, 0, 0, 40, India, >50K
23, Private, 122272, Bachelors, 13, Never-married, Adm-clerical, Own-child, White, Female, 0, 0, 30, United-States, <=50K

Bias & Fairness in NLP

A representative task – occupation prediction from biographies:

[She/He?] graduated from Lehigh University, with honours in 1998.

[Nancy/Adam?] has years of experience in weight loss surgery, patient support, education, and diabetes.



Nurse

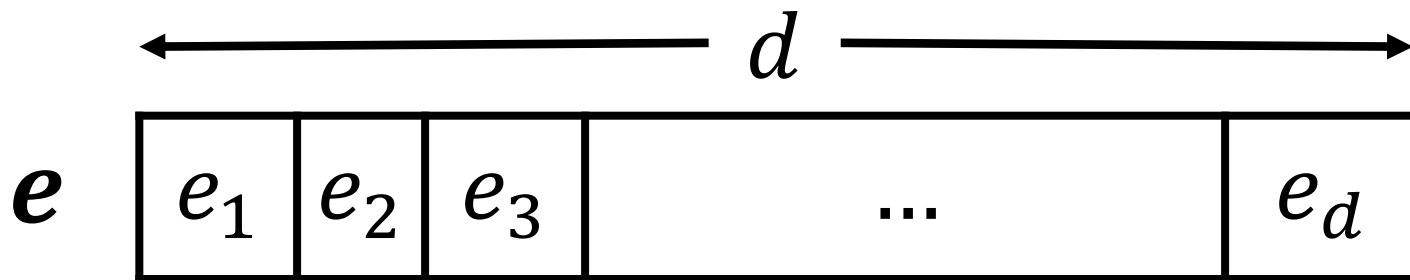
Language is inherently intertwined with
semantics and implicit meanings

Agenda

- Bias & fairness in NLP ... what? why?
- **Observing biases**
- Fairness in biography classification

Representation learning and bias

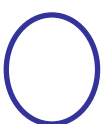
Representation learning encodes information but also may encode the **underlying biases** in data!



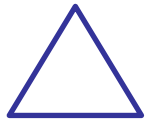
Recap

Märzen

Tesgüino

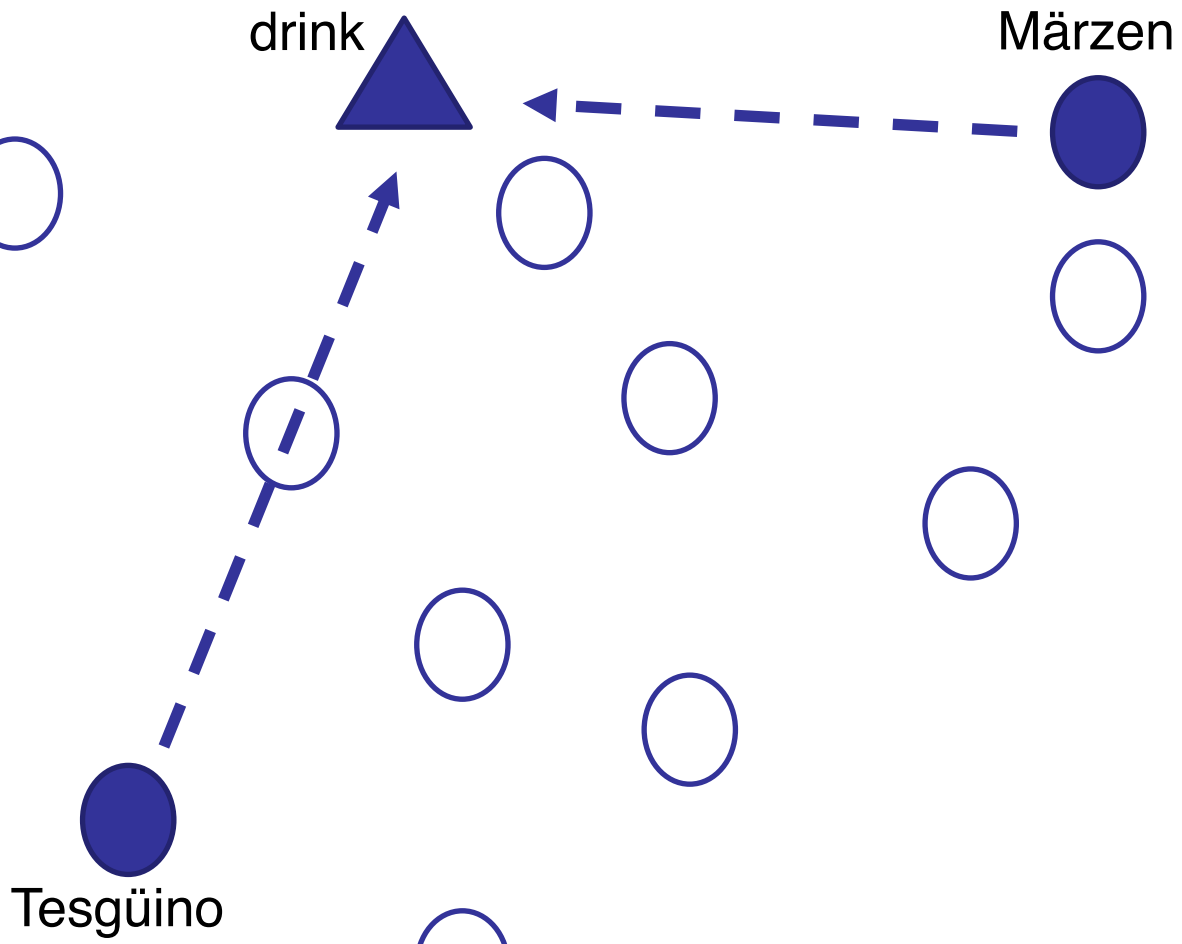


Embedding vector



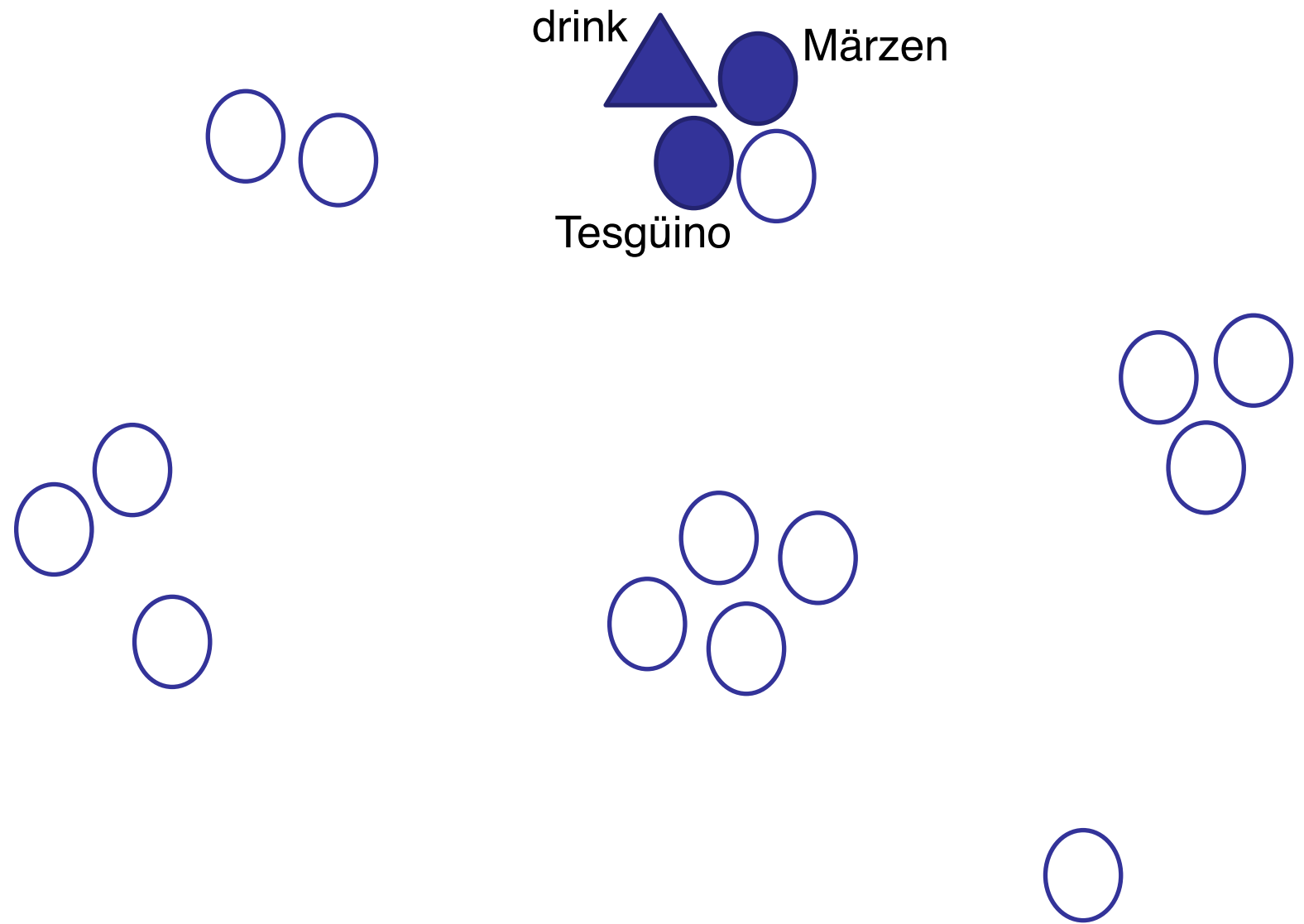
Decoding vector

Recap

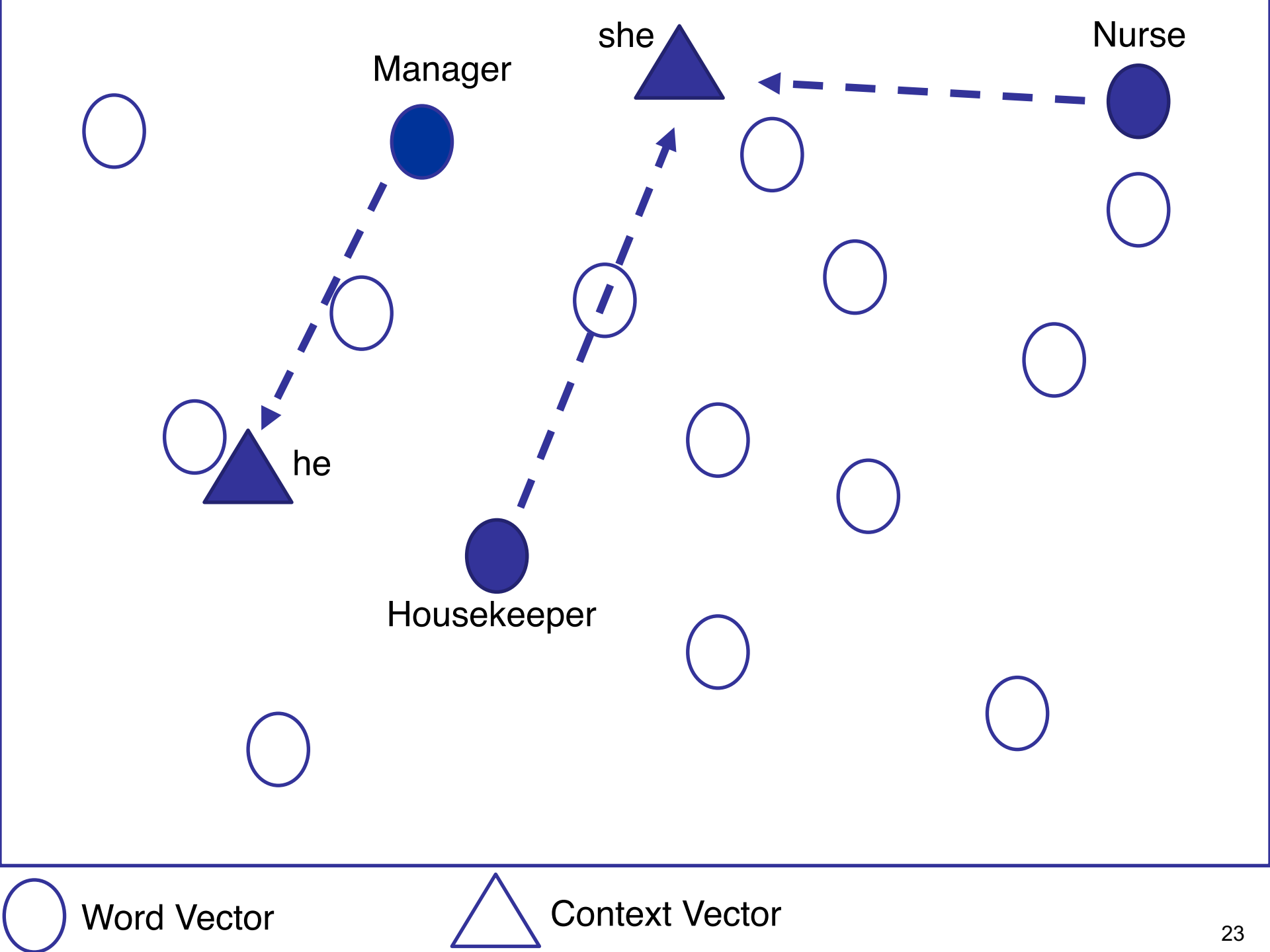


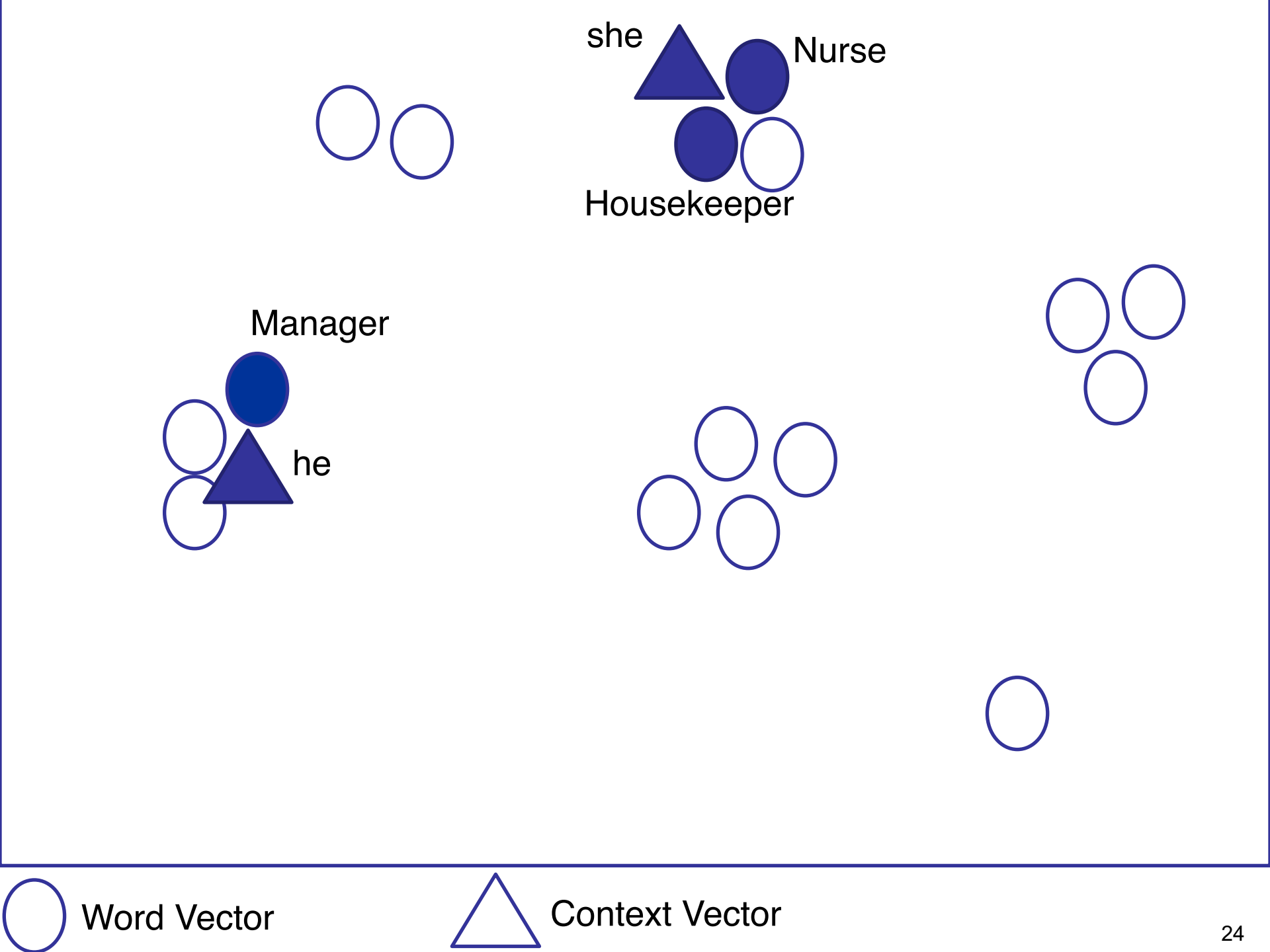
 Embedding vector  Decoding vector

Recap



 Embedding vector  Decoding vector





Bias in word analogies

- Recap – word analogy: *man* to *woman* is like *king* to ? (*queen*)

$$\mathbf{x}_{\text{king}} - \mathbf{x}_{\text{man}} + \mathbf{x}_{\text{woman}} = \mathbf{x}^*$$

$$\mathbf{x}^* \approx \mathbf{x}_{\text{queen}}$$

- Gender bias is reflected in word analogies

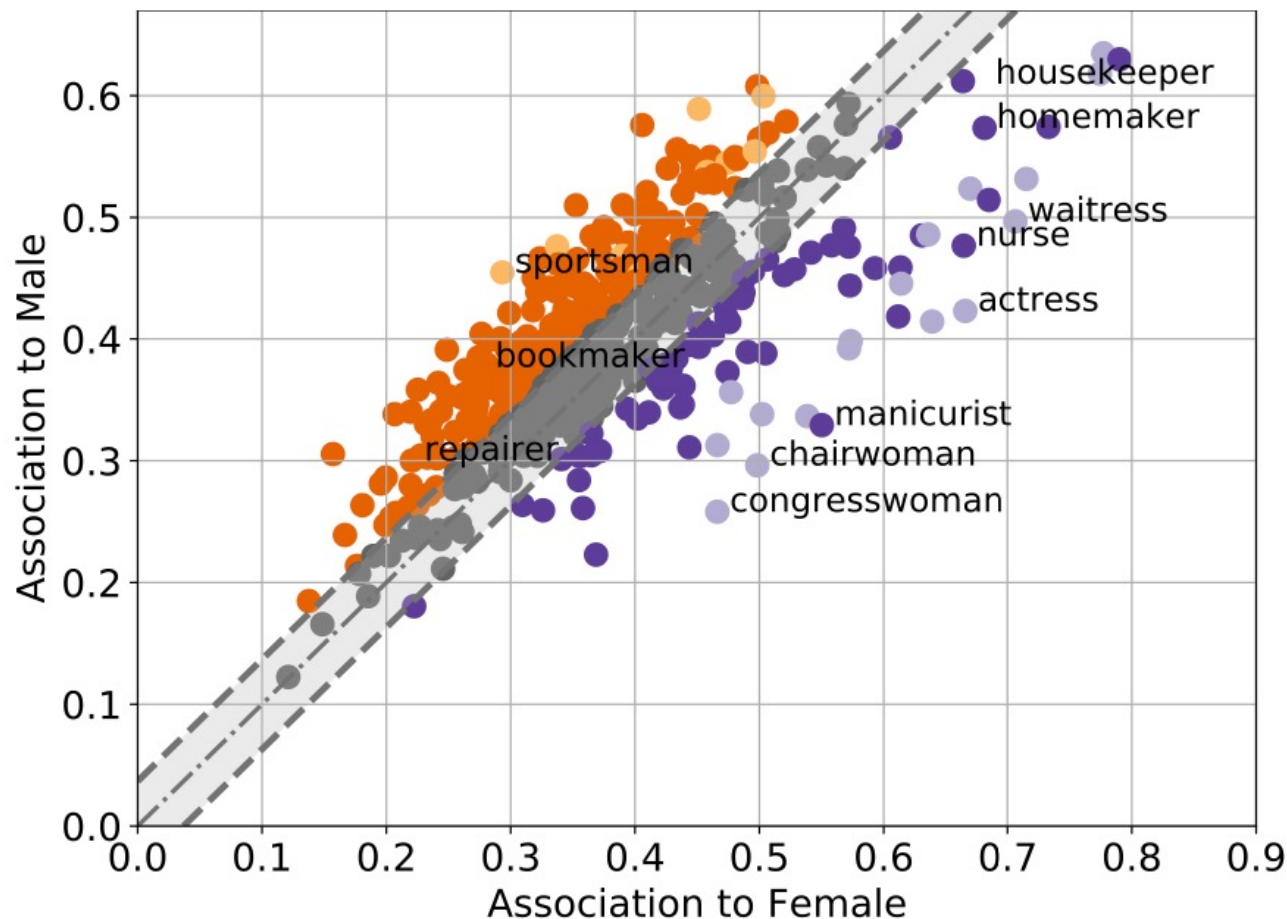
Gender stereotype *she-he* analogies

sewing-carpentry	registered nurse-physician	housewife-shopkeeper
nurse-surgeon	interior designer-architect	softball-baseball
blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
giggle-chuckle	vocalist-guitarist	petite-lanky
sassy-snappy	diva-superstar	charming-affable
volleyball-football	cupcakes-pizzas	lovely-brilliant

Gender appropriate *she-he* analogies

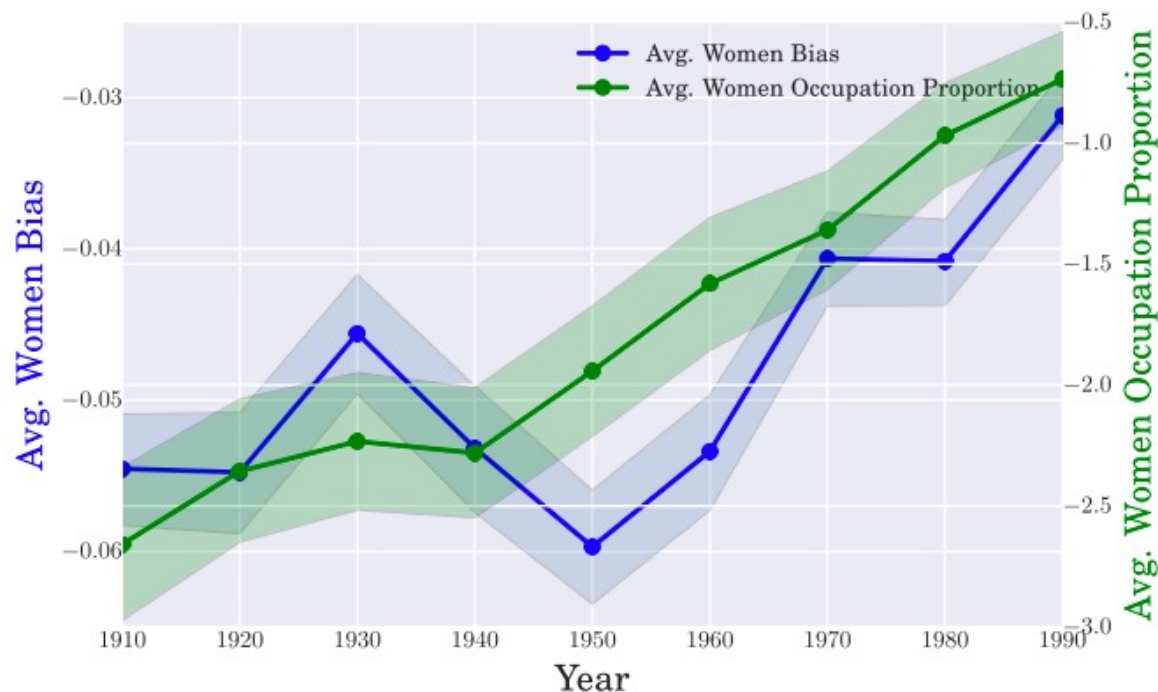
queen-king	sister-brother	mother-father
waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

Gender bias of words in a word embedding model



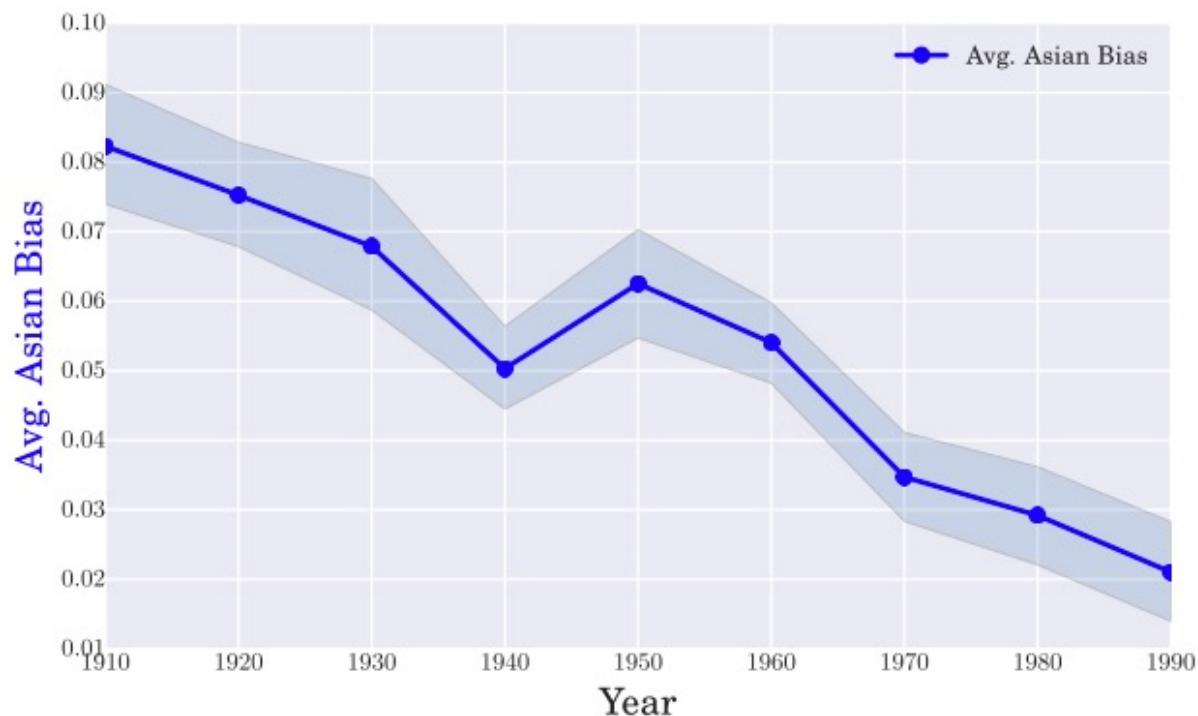
A word2vec model trained on a recent Wikipedia corpus

Word embeddings reflect societal changes!



(b) Average gender bias score over time in COHA embeddings in occupations vs the average log proportion. In blue is relative women bias in the embeddings, and in green is the average log proportion of women in the same occupations.

Word embeddings reflect societal changes!



(c) Asian bias score over time for words related to the outsiders in COHA data.

Measuring bias of a word using word embeddings

High-order method

- Bias: discrepancy of the relations of a word w (like nurse) towards two concepts \mathbb{V} and $\tilde{\mathbb{V}}$ (like female and male)
- \mathbb{V} and $\tilde{\mathbb{V}}$ are commonly defined by sets of **representative words**. For example, in a *binary* setting of gender bias:

$$\mathbb{V} = \{\textit{she}, \textit{her}, \textit{woman}, \textit{girl}, \dots\}$$

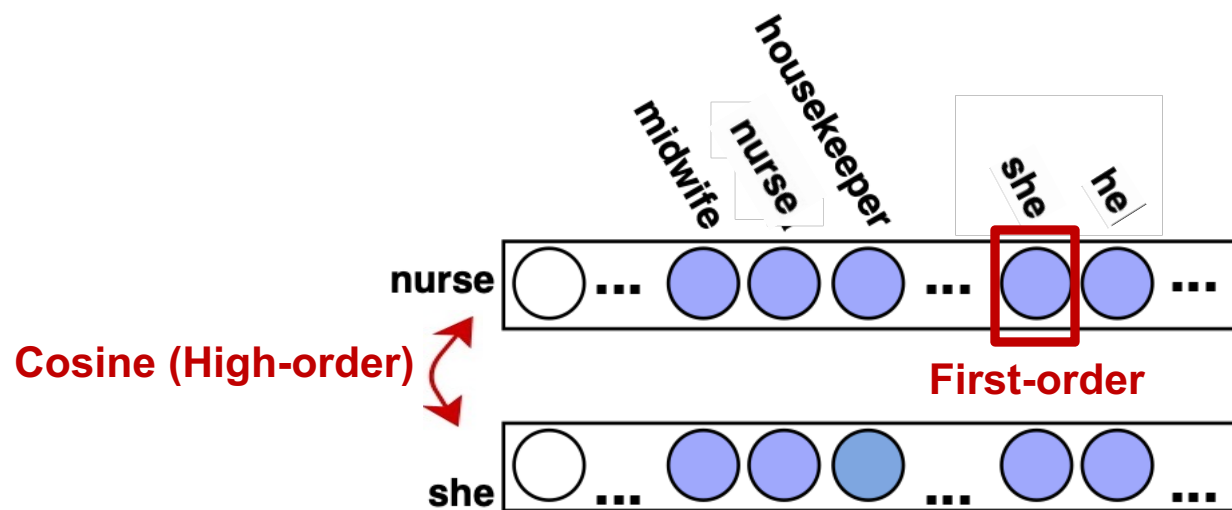
$$\tilde{\mathbb{V}} = \{\textit{he}, \textit{him}, \textit{man}, \textit{boy}, \dots\}$$

- High-order bias measurement:

$$\text{BIAS}_{\text{High}}(w) = \frac{1}{|\mathbb{V}|} \sum_{v \in \mathbb{V}} \cos(\mathbf{e}_v, \mathbf{e}_w) - \frac{1}{|\tilde{\mathbb{V}}|} \sum_{\tilde{v} \in \tilde{\mathbb{V}}} \cos(\mathbf{e}_{\tilde{v}}, \mathbf{e}_w)$$

First-order Bias Measurement

$$\text{BIAS}_{2\text{ND}}(w) = \frac{1}{|\mathbb{V}|} \sum_{v \in \mathbb{V}} \cos(e_v, e_w) - \frac{1}{|\tilde{\mathbb{V}}|} \sum_{\tilde{v} \in \tilde{\mathbb{V}}} \cos(e_{\tilde{v}}, e_w)$$



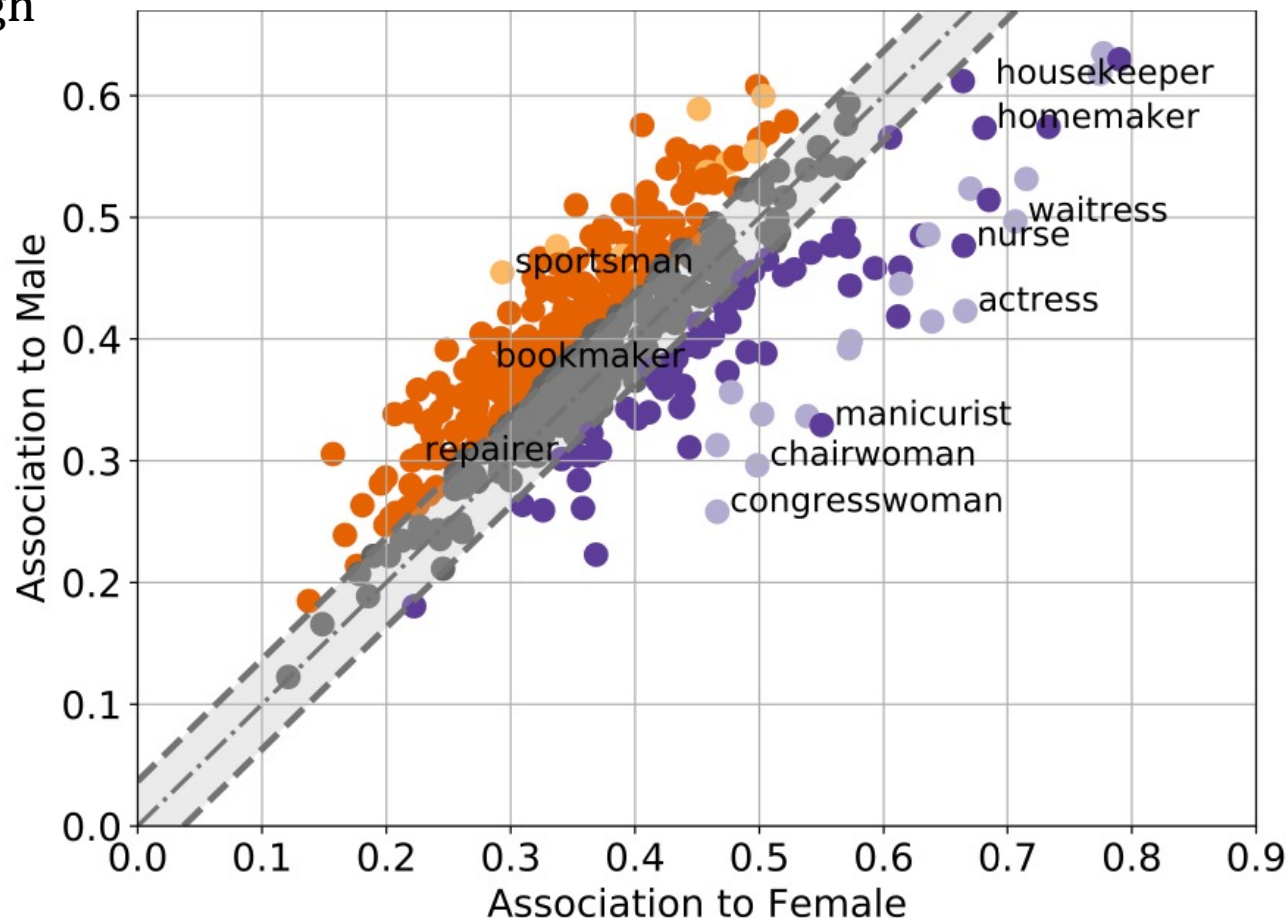
- First-order bias measurement for word w :

$$\text{BIAS}_{\text{First}}(w) = \frac{1}{|\mathbb{V}|} \sum_{v \in \mathbb{V}} f(v, w) - \frac{1}{|\tilde{\mathbb{V}}|} \sum_{\tilde{v} \in \tilde{\mathbb{V}}} f(\tilde{v}, w)$$

f provides a smoothed first-order relation between the words, achieved by reconstructing co-occurrence matrix from word embeddings and context-word embeddings

Measuring bias in WE with high-order method

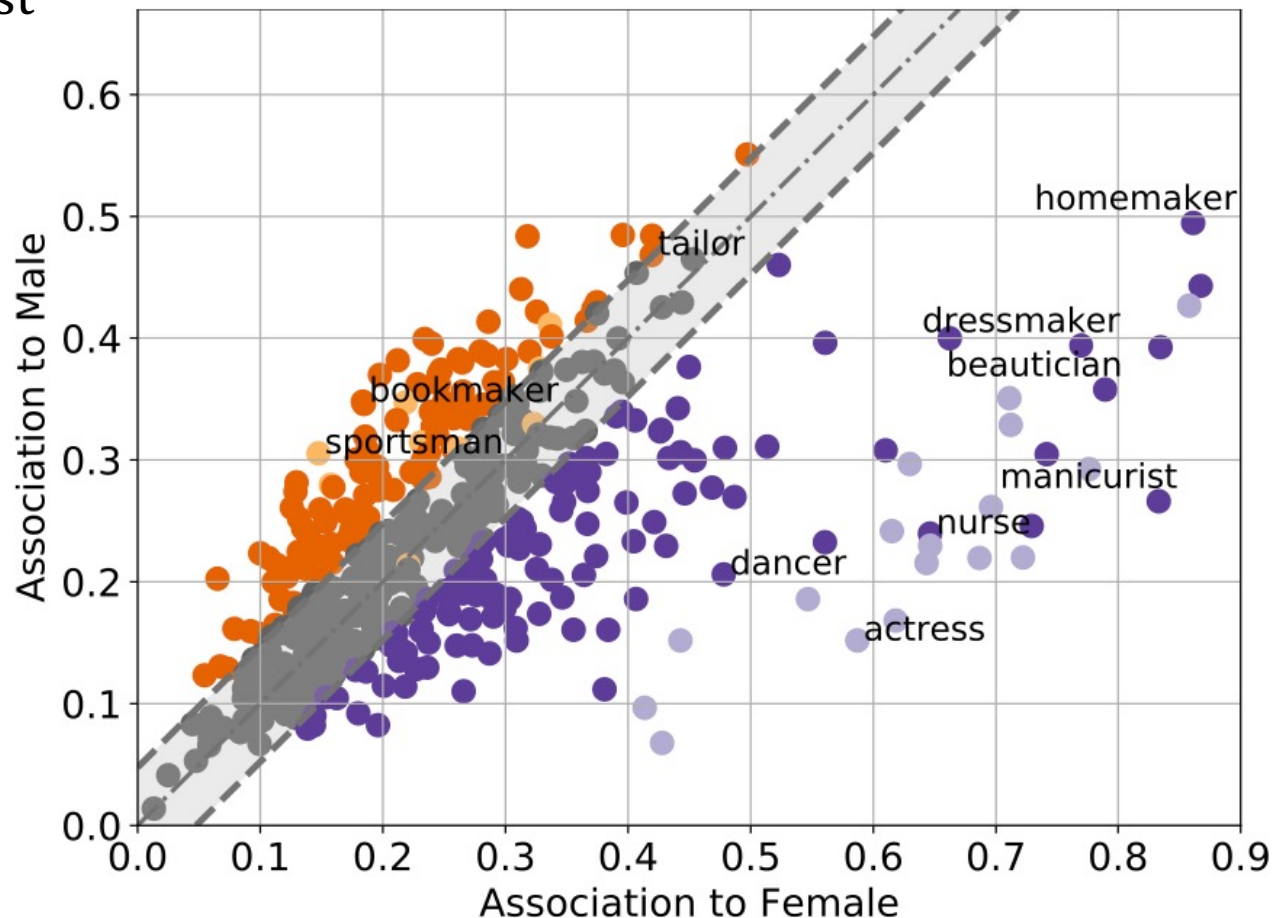
BIAS_{High}



A word2vec model trained on a recent Wikipedia corpus

Measuring bias in WE with first-order method

BIAS_{First}



A word2vec model trained on a recent Wikipedia corpus

Correlations with job market statistics

Order	Representation	Method	Labor Data		Census Data	
			Spearman ρ	Pearson's r	Spearman ρ	Pearson's r
High-Order	PMI	DIRECTIONAL	0.28	0.07	0.18	0.02
		CENTROID	0.14	0.21	0.35	0.40
		AVERAGE _{HIGH}	0.33	0.24	0.27	0.19
	PMI-SVD	DIRECTIONAL	0.05	0.07	0.00	0.00
		CENTROID	0.41	0.47	0.46	0.53
		AVERAGE _{HIGH}	0.41	0.49	0.49	0.56
→ First-Order	PMI	AVERAGE _{FIRST}	0.53	0.51	0.57	0.62
High-Order	PPMI	DIRECTIONAL	0.45	0.49	0.39	0.47
		CENTROID	0.43	0.46	0.45	0.50
		AVERAGE _{HIGH}	0.43	0.46	0.45	0.52
	PPMI-SVD	DIRECTIONAL	0.05	0.07	0.00	0.00
		CENTROID	0.41	0.47	0.46	0.53
		AVERAGE _{HIGH}	0.41	0.49	0.49	0.56
→ First-Order	PPMI	AVERAGE _{FIRST}	0.59	0.58	0.64	0.64
High-Order	SPPMI	DIRECTIONAL	0.26	0.37	0.26	0.28
		CENTROID	0.39	0.45	0.45	0.48
		AVERAGE _{HIGH}	0.32	0.40	0.44	0.48
	SPPMI-SVD	DIRECTIONAL	0.17	0.29	0.11	0.03
		CENTROID	0.28	0.35	0.39	0.43
		AVERAGE _{HIGH}	0.26	0.38	0.36	0.46
→ First-Order	SPPMI	AVERAGE _{FIRST}	0.57	0.49	0.52	0.48
High-Order	GloVe	DIRECTIONAL	0.53	0.56	0.34	0.46
		CENTROID	0.58	0.60	0.39	0.51
		AVERAGE _{HIGH}	0.60	0.60	0.39	0.51
→ First-Order	initGlove eGloVe	AVERAGE _{FIRST}	0.38	0.42	0.40	0.51
High-Order	SG	DIRECTIONAL	0.50	0.54	0.58	0.64
		CENTROID	0.55	0.57	0.60	0.65
		AVERAGE _{HIGH}	0.55	0.57	0.59	0.65
→ First-Order	eSG	AVERAGE _{FIRST}	0.66	0.61	0.67	0.70

Correlation results of the gender bias values (calculated with word embeddings) to the statistics of the portion of women in occupations

Summary

- Word embeddings capture and **encode societal biases**, reflected in the underlying corpora
 - These biases also exist in contextualized word embeddings
- Word embeddings enable the study of **societal phenomena**
 - e.g. monitoring how gender/ethnicity/etc. is perceived during time

Agenda

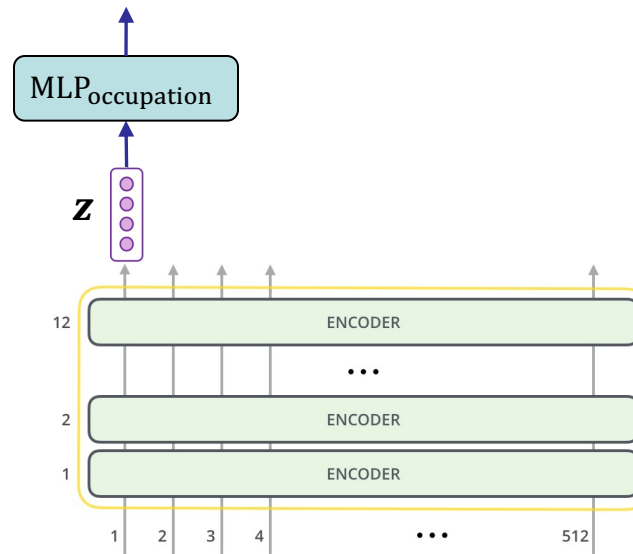
- Bias & fairness in NLP ... what? why?
- Observing biases
- **Fairness in biography classification**

Biography classification

- Predicting the occupation of a person from his/her biography
 - Explicit indications (names, pronouns) of genders are removed in the biographies

X graduated from Lehigh University, with honors in 1998. has years of experience in weight loss surgery, patient support, education, and diabetes. \longrightarrow $y = \text{Nurse}$

$y = \text{occupation}$



MLP: Multi-layer Perceptron

- Usually only one linear transformation

Data & Evaluation

Gender (sensitive attribute)	Input	Label	Prediction of a model
Male	X_1	Surgeon	Surgeon
Male	X_2	Surgeon	Surgeon
Male	X_3	Surgeon	Surgeon
Male	X_4	Surgeon	Surgeon
Male	X_5	Surgeon	Nurse
Male	X_6	Surgeon	Surgeon
Male	X_7	Surgeon	Surgeon
Male	X_8	Surgeon	Surgeon
Female	X_9	Surgeon	Nurse
Female	X_{10}	Surgeon	Surgeon
Female	X_{11}	Surgeon	Surgeon
Female	X_{12}	Surgeon	Surgeon
Male	X_{13}	Nurse	Surgeon
Male	X_{14}	Nurse	Nurse
Female	X_{15}	Nurse	Nurse
Female	X_{16}	Nurse	Nurse
Female	X_{17}	Nurse	Nurse
Female	X_{18}	Nurse	Nurse
Female	X_{19}	Nurse	Surgeon
Female	X_{20}	Nurse	Nurse

- Evaluation metric: True Positive Rate (TPR)
- TPR per occupation:

$$TPR_{occ} = \frac{\text{\# of correct Occupation}}{\text{\# of Occupation}}$$

$$TPR_{Surgeon} = \frac{10}{12} = \frac{5}{6}$$

$$TPR_{Nurse} = \frac{6}{8} = \frac{3}{4}$$

- TPR per occupation and gender:

$$TPR_{occ,gender} = \frac{\text{\# of correct for Occupation and Gender}}{\text{\# of Occupation and Gender}}$$

$$TPR_{Surgeon,Male} = \frac{7}{8}$$

$$TPR_{Surgeon,Female} = \frac{3}{4}$$

$$TPR_{Nurse,Male} = \frac{1}{2}$$

$$TPR_{Nurse,Female} = \frac{5}{6}$$

Fairness as equality in the quality of service (one possible definition)

Gender (sensitive attribute)	Input	Label	Prediction of a model
Male	X_1	Surgeon	Surgeon
Male	X_2	Surgeon	Surgeon
Male	X_3	Surgeon	Surgeon
Male	X_4	Surgeon	Surgeon
Male	X_5	Surgeon	Nurse
Male	X_6	Surgeon	Surgeon
Male	X_7	Surgeon	Surgeon
Male	X_8	Surgeon	Surgeon
Female	X_9	Surgeon	Nurse
Female	X_{10}	Surgeon	Surgeon
Female	X_{11}	Surgeon	Surgeon
Female	X_{12}	Surgeon	Surgeon
Male	X_{13}	Nurse	Surgeon
Male	X_{14}	Nurse	Nurse
Female	X_{15}	Nurse	Nurse
Female	X_{16}	Nurse	Nurse
Female	X_{17}	Nurse	Nurse
Female	X_{18}	Nurse	Nurse
Female	X_{19}	Nurse	Surgeon
Female	X_{20}	Nurse	Nurse

- A system is fair (regarding a sensitive attribute), if it provides an equal quality of service to the underlying social groups

One metric of *unfairness*:

$$\text{Unfairness}_{\text{occ}} = \text{TPR}_{\text{occ, Male}} - \text{TPR}_{\text{occ, Female}}$$

Example:

$$\text{TPR}_{\text{Surgeon, Male}} = \frac{7}{8}$$

$$\text{TPR}_{\text{Nurse, Male}} = \frac{1}{2}$$

$$\text{TPR}_{\text{Surgeon, Female}} = \frac{3}{4}$$

$$\text{TPR}_{\text{Nurse, Female}} = \frac{5}{6}$$

$$\text{Unfairness}_{\text{Surgeon}} = \frac{7}{8} - \frac{3}{4} = \frac{1}{8}$$

Unfair towards female

$$\text{Unfairness}_{\text{Nurse}} = \frac{1}{2} - \frac{5}{6} = -\frac{1}{3}$$

Unfair towards male

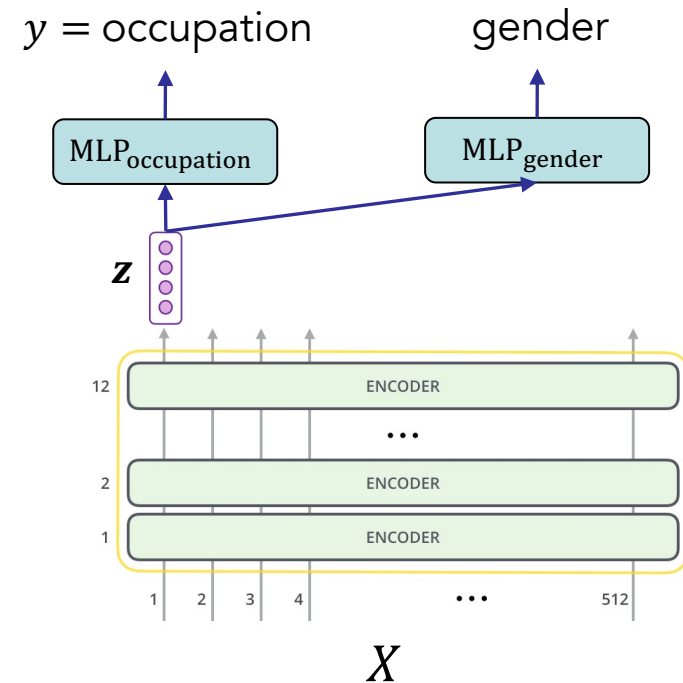
$$\text{Unfairness}_{\text{system}} = \left| -\frac{1}{8} \right| + \left| \frac{1}{3} \right|$$

Fairness as “blindness” (another possible definition)

- A system is fair (regarding a sensitive attribute), if its decisions (predictions) are agnostic to the underlying social groups
 - A fair system has no knowledge (is blind) in its encoded representations regarding the sensitive attribute

One practical approach:

- If a separate network/head, defined over \mathbf{z} , can not predict the protected attribute, the model is fair
 - In the case of binary gender, if the prediction accuracy is random (50%)



Approaching bias/unfairness

- Common approaches to reduce bias and increase fairness
 - Data curation
 - Algorithmic approaches to bias mitigation / fairness support

Algorithmic approaches:

- Pre-processing:
 - Changing/Manipulating dataset
- In-processing:
 - Impose fairness criteria to models' learning processes
 - Supporting the performance of models for minority groups
 - Removing protected information in learned embeddings
 - ...
- Post-processing
 - Changing/Rearranging model's outputs