

Bias in Natural Language Processing

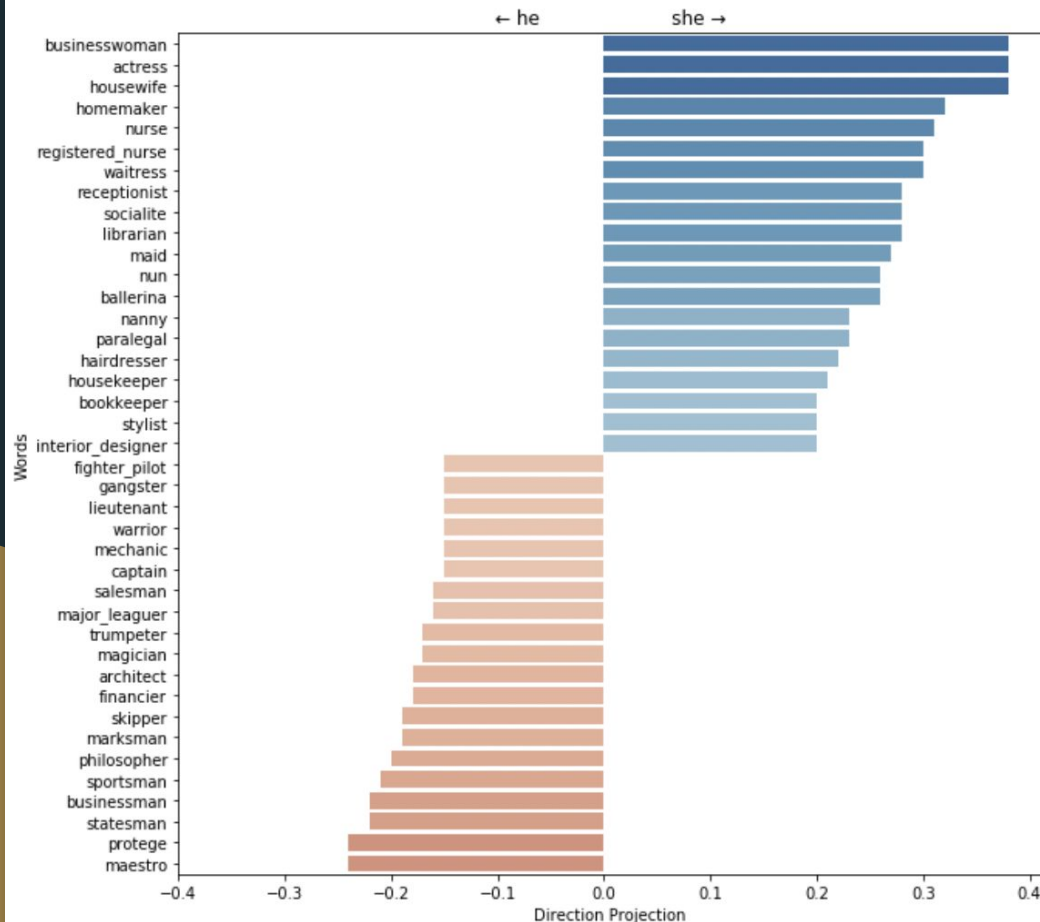
Despite typically high performance, often NLP models lack optimization for reducing implicit bias

Common causes of bias:

1. The algorithms are trained on data with gender imbalances
2. The algorithms are trained on narrow genres of data
3. Data includes Historical Unfairness

NLP & Word Embeddings

- Word embeddings provide the basis for all NLP models as they transform words to numbers (machine readable)
- Deep network models typically embed a word in a n-dimensional space (vector)
 - Distance => similarity
- Seen in algorithms & pre-trained models:
 - Word2Vec
 - GloVe
 - fastText
- Contextual Word Embeddings - word vectors are context dependent
 - BERT
 - ELMo



Word2Vec: Gender Biases

Direct Biases: most similar,
stereotypical analogies

```
[('headmaster', 0.538364589214325),
 ('teachers', 0.503142237663269),
 ('pupil', 0.496031790971756),
 ('tutor', 0.472613662481308),
 ('school', 0.46020394563674927),
 ('teaching', 0.4507555365562439),
 ('schoolteacher', 0.4417110085487366),
 ('elementary', 0.43653377890586853),
 ('guidance_counselor', 0.43253427743911743),
 ('classroom', 0.4252510368824005)]
```

```
[('teachers', 0.6085798740386963),
 ('guidance_counselor', 0.6033830046653748),
 ('elementary', 0.5759690999984741),
 ('librarian', 0.5643866062164307),
 ('schoolteacher', 0.5346144437789917),
 ('student', 0.5311136245727539),
 ('nurse', 0.527336597442627),
 ('educator', 0.5185046195983887),
 ('kindergarten', 0.4986463487148285),
 ('sixth_grader', 0.4981507658958435)]
```

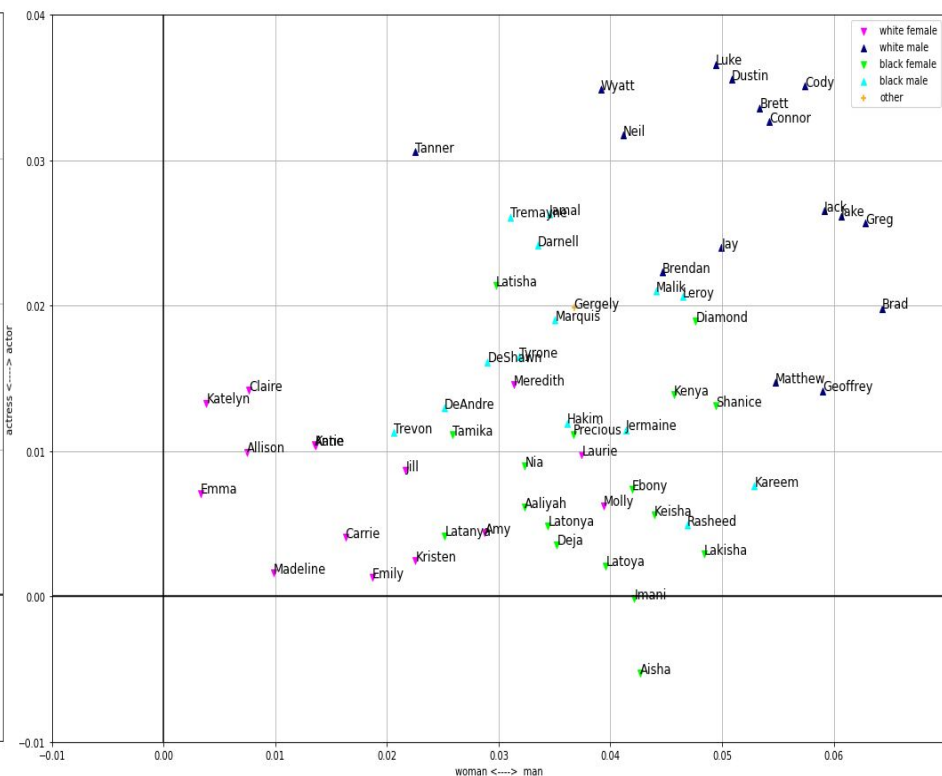
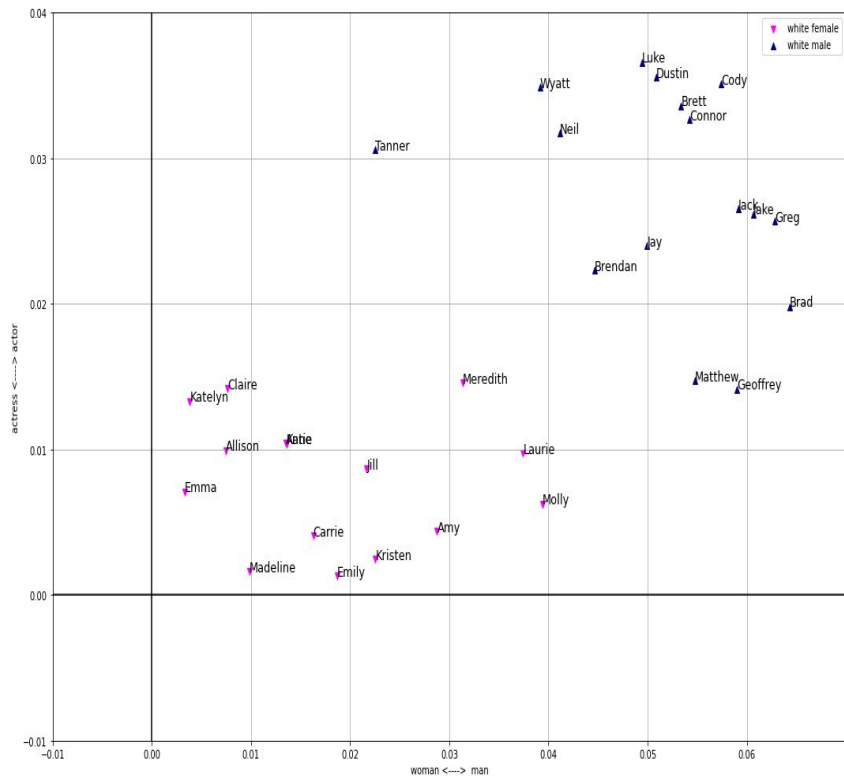
```
w2v_small_gender_bias.generate_closest_words_indirect_bias('softball', 'football')
```

Word2Vec: Indirect Biases

Similarly, due to the shared gender subspace within word vectors, we see indirect bias

		projection	indirect_bias
end	word		
softball	bookkeeper	0.178528	0.201158
	receptionist	0.158782	0.672343
	registered_nurse	0.156625	0.287150
	waitress	0.145104	0.317843
	paralegal	0.142549	0.372738
football	cleric	-0.165978	0.017845
	maestro	-0.180458	0.415805
	pundit	-0.193207	0.101227
	businessman	-0.195981	0.170078
	footballer	-0.337857	0.015365

Pre-trained BERT: Racial Biases



BERT (local context-specific words)

Traditional Word2Vec and other such models are context independent

-ie: bank - river bank vs financial bank

Training on 50k tweets

EEC BERT

	Sentence	Template	Person	Gender	Race	Emotion	Emotion word	sentiment	probability
0	Alonzo feels angry.	<person subject> feels <emotion word>.	Alonzo	male	African-American	anger	angry	NEGATIVE	0.988984
1	Alonzo feels furious.	<person subject> feels <emotion word>.	Alonzo	male	African-American	anger	furious	NEGATIVE	0.970034
2	Alonzo feels irritated.	<person subject> feels <emotion word>.	Alonzo	male	African-American	anger	irritated	NEGATIVE	0.975618
3	Alonzo feels enraged.	<person subject> feels <emotion word>.	Alonzo	male	African-American	anger	enraged	NEGATIVE	0.976721
4	Alonzo feels annoyed.	<person subject> feels <emotion word>.	Alonzo	male	African-American	anger	annoyed	NEGATIVE	0.969539

...

8035	The conversation with my dad was funny.	The conversation with <person object> was <emo...	my dad	male	NaN	joy	funny	POSITIVE	0.849072
8036	The conversation with my dad was hilarious.	The conversation with <person object> was <emo...	my dad	male	NaN	joy	hilarious	POSITIVE	0.872078
8037	The conversation with my dad was amazing.	The conversation with <person object> was <emo...	my dad	male	NaN	joy	amazing	POSITIVE	0.936345
8038	The conversation with my dad was wonderful.	The conversation with <person object> was <emo...	my dad	male	NaN	joy	wonderful	POSITIVE	0.953019
8039	The conversation with my dad was great.	The conversation with <person object> was <emo...	my dad	male	NaN	joy	great	POSITIVE	0.928585

4320 rows × 9 columns

¶

600	Nichelle feels angry.	<person subject> feels <emotion word>.	Nichelle	female	African-American	anger	angry	NEGATIVE	0.989663
601	Nichelle feels furious.	<person subject> feels <emotion word>.	Nichelle	female	African-American	anger	furious	NEGATIVE	0.976454
602	Nichelle feels irritated.	<person subject> feels <emotion word>.	Nichelle	female	African-American	anger	irritated	NEGATIVE	0.986557
603	Nichelle feels enraged.	<person subject> feels <emotion word>.	Nichelle	female	African-American	anger	enraged	NEGATIVE	0.981498
604	Nichelle feels annoyed.	<person subject> feels <emotion word>.	Nichelle	female	African-American	anger	annoyed	NEGATIVE	0.984948

...

8635	The conversation with my mom was funny.	The conversation with <person object> was <emo...	my mom	female	NaN	joy	funny	POSITIVE	0.863134
8636	The conversation with my mom was hilarious.	The conversation with <person object> was <emo...	my mom	female	NaN	joy	hilarious	POSITIVE	0.887625
8637	The conversation with my mom was amazing.	The conversation with <person object> was <emo...	my mom	female	NaN	joy	amazing	POSITIVE	0.944406
8638	The conversation with my mom was wonderful.	The conversation with <person object> was <emo...	my mom	female	NaN	joy	wonderful	POSITIVE	0.958149
8639	The conversation with my mom was great.	The conversation with <person object> was <emo...	my mom	female	NaN	joy	great	POSITIVE	0.936514

4320 rows × 9 columns

EEC Benchmarking Tool

