



Gender Bias in Natural Language Processing Across Human Languages

Abigail Matthews¹, Isabella Grasso², Christopher Mahoney², Yan Chen², Esmá Wali², Thomas Middleton², Mariama Njie³, Jeanna Neefe Matthews²,
¹University of Wisconsin-Madison, ²Clarkson University, ³Iona College

Abstract

Natural Language Processing (NLP) systems are at the heart of many critical automated decision-making systems making crucial recommendations about our future world. Gender bias in NLP has been well studied in English, but has been less studied in other languages. In this paper, a team including speakers of 9 languages – Chinese, Spanish, English, Arabic, German, French, Farsi, Urdu, and Wolof – reports and analyzes measurements of gender bias in the Wikipedia corpora for these 9 languages. We develop extensions to profession-level and corpus-level gender bias metric calculations originally designed for English and apply them to 8 other languages, including languages that have grammatically gendered nouns including different feminine, masculine, and neuter profession words. We discuss future work that would benefit immensely from a computational linguistics perspective.

Background/Motivation

- Corpora of human language are regularly fed into machine learning systems as a key way to learn about the world.
- NLP plays a significant role in speech recognition, text translation, and autocomplete.
- NLP is the heart of many critical automated decision systems making crucial recommendations about our future world.
- Systems are taught to identify spam email, suggest medical articles or diagnoses related to a patient's symptoms, sort resumes based on relevance for a given position
- Key component of critical decision making systems in areas such as criminal justice, credit, housing, allocation of public resources and more.
- Expanding work on quantify gender bias done only in English to other languages

Future Work/ Opportunities for Collaboration

- Experiments with different defining sets, both specific to one language and the best set across languages
- Experiments with changes in defining sets and profession sets over time
- Approaches to dealing with disambiguation of terms
- Identifying different corpora beyond Wikipedia. Even within one language, it would be interesting to examine collections with different emphasis such as gender of author, different time periods, different genres of text, different country of origin, etc.

Defining Set

English	Chinese	Spanish	Arabic	German	French	Farsi	Urdu	Wolof
woman	女人	mujer	المرأة	Frau	femme	زنان	عورت	Jigéen
man	男人	hombre	الرجل	Mann	homme	مرد	مرد	Goor
daughter	女儿	hija	ابنة	Tochter	filie	دختر	دختر	Doom ju gior
son	儿子	hijo	ولد	Sohn	filis	پسر	پسر	Doom ju gior
nephew	侄子	nieto	ابن عم	Königin	reine	ملکه	ملکه	Jahur baar
father	父亲	padre	أب	Vater	père	پدر	پدر	Baay
girl	女孩	nina	ابنة	Mädchen	filie	دختر	دختر	Jaxna
boy	男孩	nino	مسي	Junge	garçon	پسر	پسر	Xale bu gior
queen	女王	reina	ملکه	Königin	reine	ملکه	ملکه	Jahur baar
king	国王	rey	ملك	König	roi	پادشاه	پادشاه	Baar
wife	妻子	esposa	زوجة	Ehefrau	épouse	همسر	همسر	Jibar
husband	丈夫	esposo	الزوج	Ehemann	marri	شوهر	شوهر	jikkar
madam	女士	señora	سيدة	Dame	madame	خانم	خانم	Nilawri
sir	男士	señor	سيدي	Herr	monsieur	آقا	آقا	Goopil

Defining Set

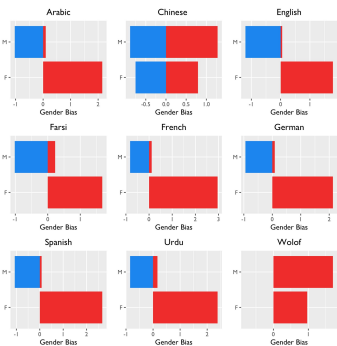
The defining set is a list of gendered word pairs used to define what a gendered relationship looks like. Bolukbasi et al's original defining set contained 10 English word pairs (she-he, daughter-son, her-his, mother-father, woman-man, gal-guy, Mary-John, girl-boy, herself- himself, and female-male). We began with this set, but made substantial changes in order to compute gender bias effectively across 9 languages. Specifically, we removed 6 of the 10 pairs, added 3 new pairs, and translated the final set into 8 additional languages.

We use Bolukbasi et al's formula for direct gender bias: where N represents the list of profession words, g represents the gender direction calculated, w represents each profession word, and c is a parameter to measure the strictness of the bias.

$$\text{DirectBias}_c = \frac{1}{|N|} \sum_{w \in N} |\cos(\vec{w}, g)|^c$$

We present the gender bias scores, calculated as described above according to Bolukbasi et al's methodology, for each of our 14 defining set words (7 pairs) across 9 languages. Female gender bias is represented as a positive number (red bar) and male gender bias is represented as a negative number (blue bar). Not all defining set words occur in the Wikipedia corpus for Wolof. In some cases, this is because they are multi-word phrases and in other cases, this is likely because of the small size of the corpora.

We aggregate the gender bias for all the male words (sir, husband, king, etc.) and all the female words (madam, wife, queen, etc.) This presentation emphasizes several key aspects of the results.



Some key observations:

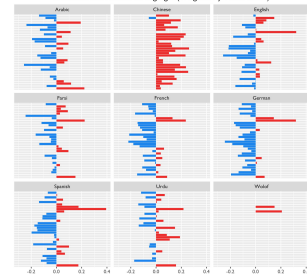
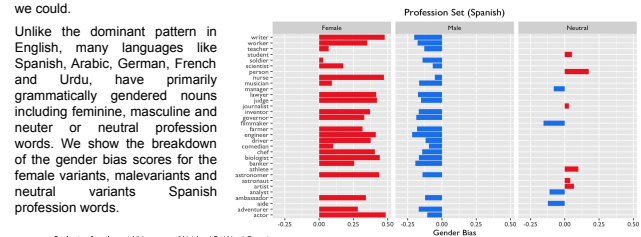
- Spanish, English, Arabic, German, French, Farsi, and Urdu, that the female words are female leaning and that most male words and male leaning as one might expect.
- Interesting exception of husband in all of these languages and also man in Farsi.
- Wolof has troubles primarily due to the small size of the corpus.
- Chinese has trouble due to difficulty finding a dominant PCA that isolates the gender direction (more info in the paper on that).

Profession Set

Profession Set

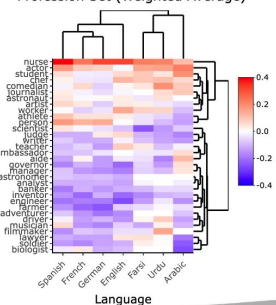
We began with Bolukbasi et al's profession word set in English, but again made substantial changes in order to compute gender bias effectively across 9 languages. Bolukbasi et al. had an original list of 327 profession words, including some words that would not technically be classified as professions like scientist or drug addict. We narrowed this list down to 32 words including: nurse, teacher, writer, engineer, scientist, manager, driver, banker, musician, artist, chef, filmmaker, judge, comedian, inventor, worker, soldier, journalist, student, athlete, actor, governor, farmer, person, lawyer, adventurer, aide, ambassador, analyst, astronaut, astronomer, and biologist. We tried to choose a diverse set of professions from creative to scientific, from high-paying to lower-paying, etc. that occurred in as many of the 9 languages as we could.

Unlike the dominant pattern in English, many languages like Spanish, Arabic, German, French and Urdu, have primarily grammatically gendered nouns including feminine, masculine and neuter or neutral profession words. We show the breakdown of the gender bias scores for the female variants, malevariants and neutral variants Spanish profession words.



We compare these profession-level gender bias scores across all 9 languages using the weighted average (weighted by word count). Notice the similarities in patterns between Spanish, English, Arabic, German, French, Farsi and Urdu. For example, nurse is often highly female biased and engineer is often highly male biased.

Profession Set (Weighted Average)



We use Pearson's correlation for cluster analysis to examine 7 languages, omitting Chinese and Wolof because of problems with PCA and corpora size respectively. This exploratory analysis provokes a number of questions for future work including: How do linguistics inform the bias outputs (e.g. If English is a mixture of Germanic and Latin languages, is that why it is clustered with those languages even though it is not gendered?) and How does a language being inherently gendered affect the resulting bias of an NLP model in that language?