

Gendered Grammar or Ingrained Bias?

Exploring Gender Bias in Icelandic Language Models

STEINUNN RUT FRÍÐRIKSDÓTTIR
HAFSTEINN EINARSSON

Biases in large language models

- Large language models are trained on a myriad of data
- More data is usually linked to better outputs
- Can result in the models absorbing all sorts of prejudices
- Sólmundsdóttir et al. (2022) showed pronounced gender bias in Icelandic machine translation models

English		Icelandic
I am strong	×	ég er sterkur
I am weak		Ég er veik
I am clever		Ég er snjall
I am stupid		ég er heimsk
I am faithful		Ég er trúr
I am unfaithful		Ég er ótrú
I am confident		ég er sjálfsöruggur
I am insecure		Ég er óörugg
I am interesting		Ég er áhugaverður
I am uninteresting		Ég er óáhugaverð

Did you mean: I am strong I am weak I a...

What is *gender bias*?

- The tendency of LMMs to generate or perpetuate gender stereotypes
- Can lead to various types of harm
 - Reinforcement of harmful societal norms
 - Dismissal of individuals that fall outside of the norms
 - Feelings of exclusion
- Can cause direct harm when used for downstream tasks



The objective

- To investigate the presence of gender bias within language models trained on Icelandic
- By focusing specifically on occupation-related terms, assess whether these models mirror the gender distributions observed in the Icelandic job market
- To achieve this, we cross-reference our findings with distribution data obtained from Statistics Iceland



Statistics Iceland

Icelandic and the generic masculine

- Icelandic is highly gendered
 - All words that inflect by case also inflect by gender
 - Doesn't reflect societal biases, just morphological agreement
- Icelandic favors the masculine when referring to a group of mixed-gendered people
 - Most occupational terms are grammatically masculine
 - 381 out of 394 in this case
 - A debate between two feminist movements
 - How does the language affect the results?



Number	Karlkyn (MASCULINE)	Kvenkyn (FEMININE)	Hvorkyn (NEUTER)
1	einn	ein	eitt
2	tveir	tvær	tvö
3	þrír	þrjár	þrjú
4	fjórir	fjórar	fjögur

Method

1

<mask> is a nurse

He: 0.15
She: 0.85
(They: 0.0)

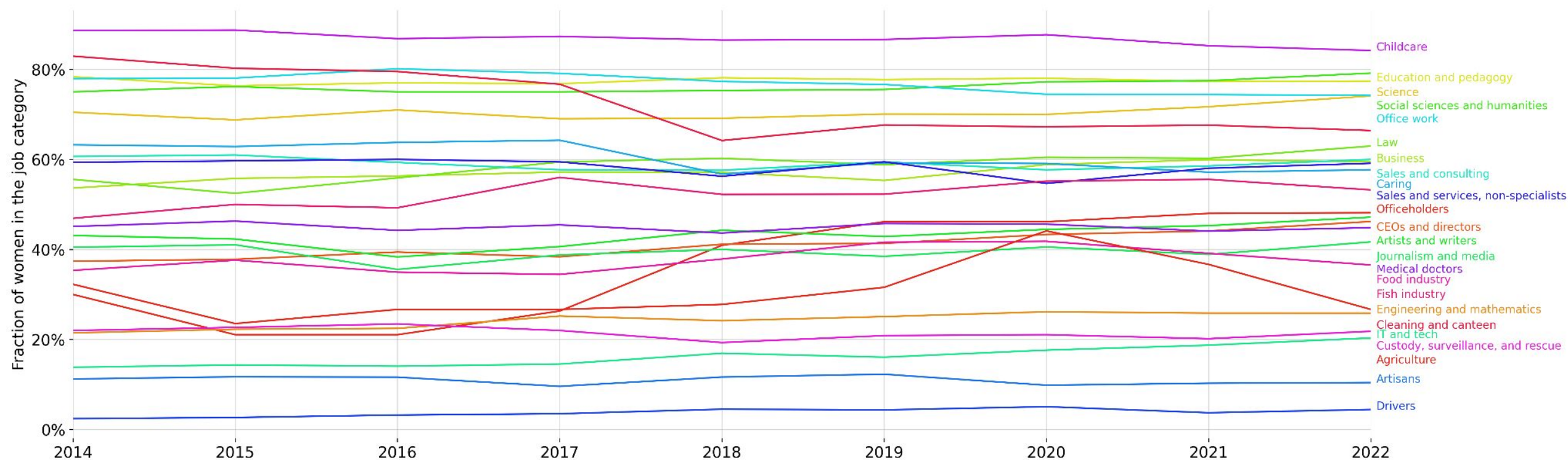
- IceBERT
- IceBERT-igc
- IceBERT-ic3
- IceBERT-xlmr-ic3
- ScandiBERT

2

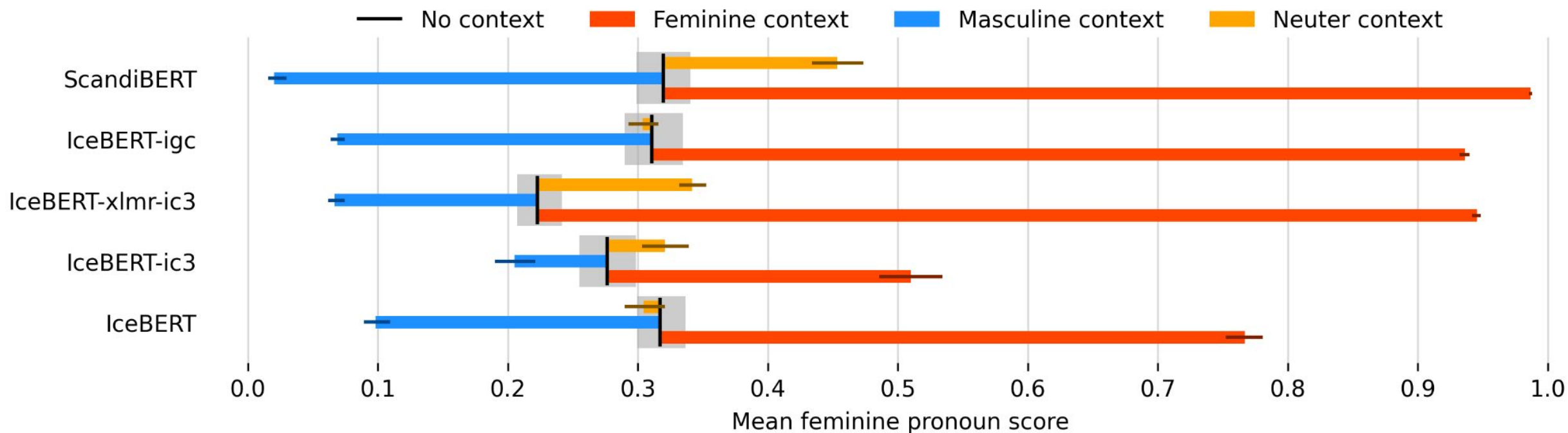
This is John. <mask> is a nurse

He: 0.35
She: 0.65
(They: 0.0)

Information from Statistics Iceland



Overall trends and the influence of context

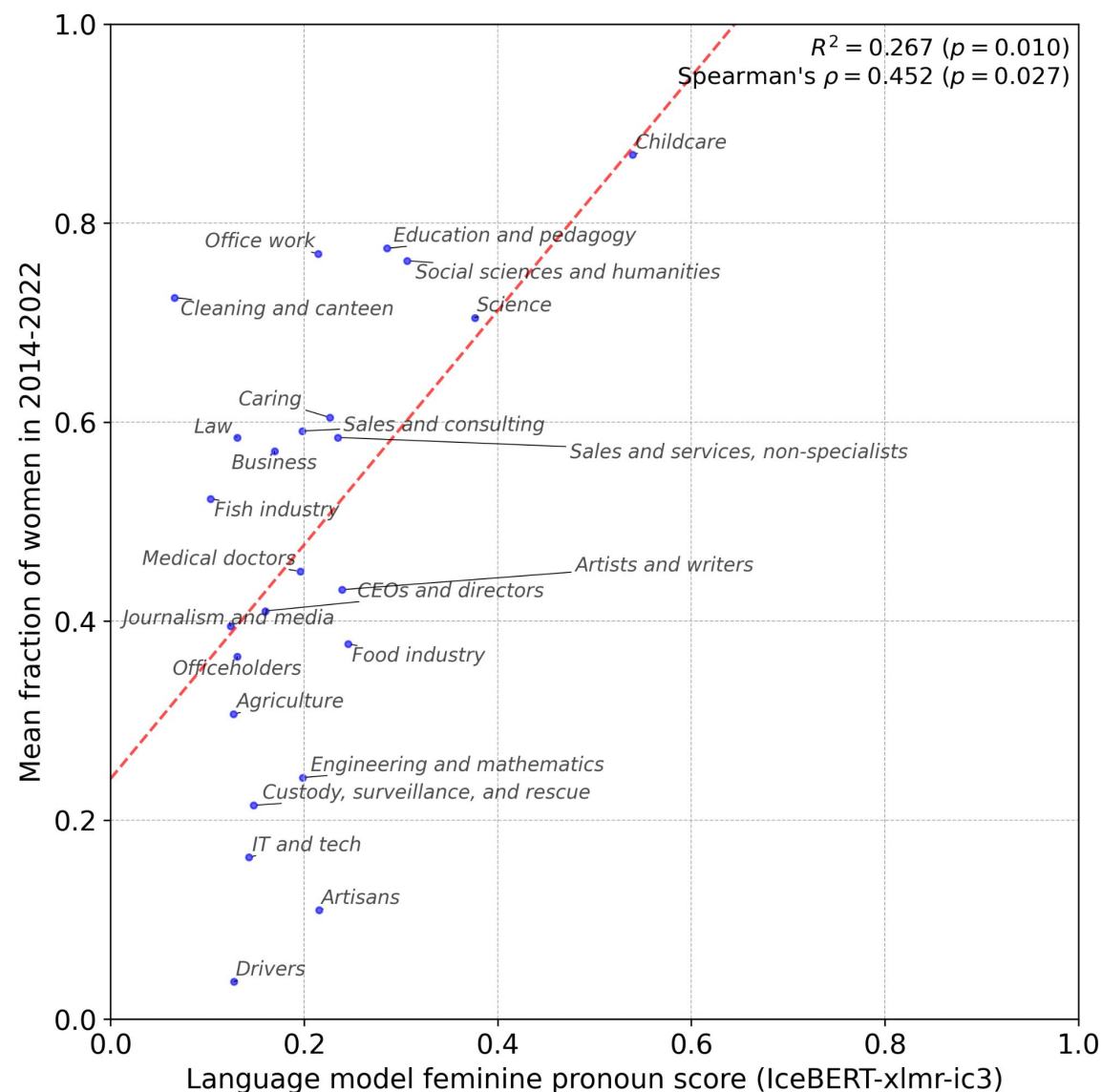


But note!

	Masc	Fem	Neut
Proper names	52,4%	30%	17,6%
Pronouns singular	29%	14,5%	56,5%
Pronouns plural	55,3%	16,6%	28,1%

Table 3: Grammatical gender distribution in the Icelandic Gigaword Corpus, used to train three of our models.

Statistical analyses



Model	r	ρ
IceBERT	0.40 (0.056)	0.36 (0.086)
IceBERT-ic3	0.46 (0.025)	0.42 (0.041)
IceB.-xlmr-ic3	0.52 (0.010)	0.45 (0.027)
IceBERT-igc	0.50 (0.012)	0.33 (0.110)
ScandiBERT	0.49 (0.016)	0.58 (0.003)

Table 2: The correlation coefficients (r), Spearman's rank correlation coefficients (ρ) and significance levels in brackets for different models when evaluating their relationship to job market data from Statistics Iceland.

Notable results

- Overall masculine bias without context
- Female-centric occupations are not overpredicted
- Compounds ending with *maður* (e. *man*), *smiður* (e. *smith*) and *meistari* (e. *master*) are more likely to be masculine
- Lower-ranking university positions are more likely to be feminine but *prófessor* (e. *full professor*) is >80% likely to be masculine
- Kindergarten, elementary and secondary school teachers are more likely to be feminine but university teachers masculine

Summary

- All models have a masculine bias but most can be mended with context
- Complex interplay of societal, linguistic and data selection biases
- We only have job categories, not individual occupations
- Would be interesting to compare our results to those of a language without a grammatical gender
- The dangers of aging data: Society evolves but the models might not be aware of it!

