# A Context-Aware System for Bias Identification in Job Advertisements using Natural Language Processing (NLP)

Vogt Luise
Guan Zhize
Nasev Veselin
Guo Yonghui

Maastricht University

# Agenda

- Introduction
  - Context and Motivation
  - Problem Statement and Research Questions
- Related Work
- Methodology
  - The Data&Data Prepcocessing
  - Data Annotation
    - Generic He/Generic She
    - Behavioral Stereotypes
    - Societal Stereotypes
    - Explicit Making of Sex
- Result & Analysis
  - Identification System
  - LIME Explaniability
- Future Work
  - Suggestion System
  - Other Bias Types
  - Modeling

# Introduction

# Context and Motivation

Bias Type:
Behavioral Stereotypes

Bias Type:
Generic He/Generic She

We are looking for a young and **driven** candidate who can bring innovation to our organization, and is a true team player for the rest within the organization. **He** needs to master relevant skills and techniques and be passionate about **his** job. We are still only 1% done at Facebook – this team is inventing every day and it takes tenacity, **bravery** and the ability to see the big opportunities to thrive.

"Independent"
is a neutral word that you can use instead of "bavery"

# Problem Statement and Research Questions

1. What kinds of biased language are commonly identified in job advertisements?
2. What words are related to the most common biases?
3. How can a context-aware natural language processing tool be used to classify different types of biases in job applications?
4. What are key predictors that explain the prediction for a particular class of bias?

**Maastricht University**

# Related Work

# Related Work



Figure 1: Thesis cover of Identifying Possible Indicators in Job Advertisments ([1])

Contributions:
- Five types of bias covering gender, rasical and other aspects
- Compared performance of different supervised machine learning classiers

Points to improve:
- The taxonomy of bias types
- Quality of annotations
- Context-aware model

# Related Work

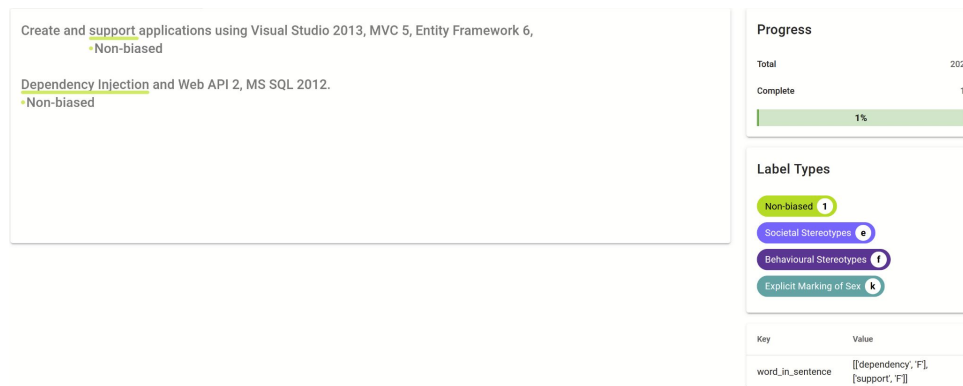**The Taxonomy of Gender Bias**

- Five parent bias types
  - Generic Pronouns
  - Stereotyping Bias
  - Sexism
  - Exclusionary Terms
  - Semantic Bias

- Based on normal English text but not job advertisements

**Maastricht University**

| Bias Type | Bias Subtype | Example | Methodology |
|---|---|---|---|
| Generic Pronouns | Generic He | A programmer must carry his laptop with him to work. | Supervised Learning |
| | Generic She | A nurse should ensure that she gets adequate rest. | Supervised Learning |
| Stereotyping Bias | Societal Stereotypes | Senators need their wives to support them throughout their campaign. | Supervised Learning |
| | Behavioural Stereotypes | The event was kid-friendly for all the mothers working in the company. | Supervised Learning |
| Sexism | Hostile Sexism | Women are incompetent at work. | Supervised Learning |
| | Benevolent Sexism | They're probably surprised at how smart you are, for a girl. | Supervised Learning |
| Exclusionary Terms | Explicit Marking of Sex | Chairman, Businessman, Manpower, Cameraman | Lexicon-Based |
| | Gendered Neologisms | Man-bread, Man-sip, Man-tini | Lexicon-Based |
| Semantic Bias | Metaphors | "Cookie": lovely woman. | Supervised Learning |
| | Old Sayings | A woman's tongue three inches long can kill a man six feet high. | Supervised Learning |

Figure 2: Overview of the taxonomy and link to detection methodology ([2])

8

# Related Work

**Data Annotation and Augmentation**

- Doccano
  - Open source text annotation tool
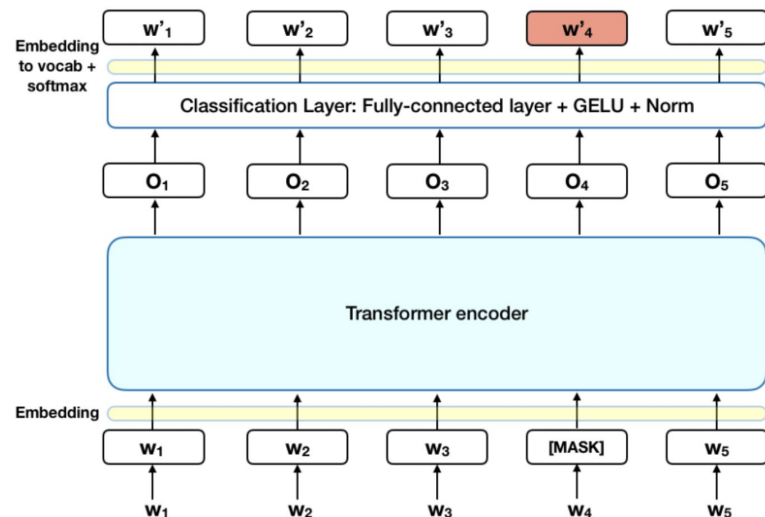  - Easy to use



- FlashText
  - Python moudle
  - Replace keywords in sentences
  - Biased word list
- Masked language model with BERT
  - Replace random words in sentences
  - Not changing the context too much

**Maastricht University**

# Related Work

**BERT model**

- Bidirectional
- Pre-trained on two tasks
- Masked language modeling
  - Hide one word
  - Predict which word has been hidden based on context
- Next sentence prediction
  - Predict whether two given sentences have a logical, sequential connection



**Maastricht University**

# Methodology

# Methodology - Data and Preprocessing

- Datasets
  - EMSCAD[3]
    - Employment Scam Aegean Dataset
    - Public available
    - 17,880 real-life job advertisements and 866 fake advertisements
  - Adzuna[4]
    - Dataset from Kaggle
    - 242,138 job adertisements

- Preprocessing
  - Remove white spaces and new rows
  - Remove HTML characters
  - Remove email addresses and phone numbers
  - Split job descriptions into sentences
  - Tokenize each sentence
  - Remove too long and too short sentences
  - Filter sentences that contain at least one biased word

**Maastricht University**

# Methodology - Generic He/She

- Found a dataset that addresses this type of gender bias
  - 1585 sentences that address (he, she, his, hers, himself, herself)

**Example:** Can you identify a programmer based on **his** code?

- Annotated 300 sentences from the job description dataset
  - Most of this sentences are not bias
  - Machine learning model can learn the difference

**Example:** John is our new programmer, **his** diligence and hard work….

**Maastricht University**

# Methodology - Behavioral Stereotypes

- Behavioral stereotypes attributes the behaviour of someone to its gender

**Example:**
Candidates must be **compassionate** and **understanding** to difficult situations while keeping focus on task at hand in a fast-paced collaborative environment.

Two methods to generate dataset:
  - Annotated 1000 sentences from the EMSCAD dataset
  - Repalce biased words on the list
  - BERT + mask

  - Manually labelled 500 sentences form Kaggle dataset
  - BERT + mask

**Maastricht University**

# Methodology - Societal Stereotypes

- Not very common in job descriptions
  - Just a few examples from Jad Doughman's paper

**Examples**:

- Senators need their wives to support them throughout their campaign.
- The event was kid-friendly for all the mothers working in the company.

- Two methods to generate more examples
  - Combine different occupations with men or women together
  - Use techniques like fill-mask and GPT-2 to replace some parts of a sentence,

**Examples:**

- Professors are men.
- Doctors are women.
- Senators need their wives to support them and the country — and not try to take away their rights.
- The program was kid-friendly for all the mothers involved in the company.

# Methodology - Explicit Marking of Sex

- mentions specific gender
- excludes other genders


- base lexicon
- word2vec
- find similar words
- manually check
- filter sentences

```
1   word
2   lady_fashion
3   gentleman
4   foreman
5   handyman
6   draughtsman
7   tradesman
8   manpower
9   chairman
10  landlord
11  manmanagement
12  acuman
13  man_management
```

**Maastricht University**

# Result & Analysis

# Results and Analysis - Dataset

**Distribution of labels**



Societal Stereotypes — 9.7% (311)
Explicit Marking of Sex — 6.2%
Generic he/she — 37.3% (1,882)
Behavioural Stereotypes — 46.9% (2,364)
488

# Results and Analysis - Modeling

"It is necessary a dedicated and ambitious food development director to join our team."

⬇

Vector of length 768 for each token

⬇

Neural network with single input layer of 768 neurons

⬇

Output layer of length number of prediction classes

⬇

Softmax function on the output layer

**Maastricht University**

# Results and Analysis - Validation

| Precision | Recall | F1 |
|:---:|:---:|:---:|
| 0.98983 | 0.98983 | 0.98983 |

| | Type | O | Behavioural Stereotypes | Explicit Marking of Sex | Societal Stereotypes | Generic he | Generic she |
|---|---|---|---|---|---|---|---|
| 0 | Precision | 0.995627 | 0.861968 | 0.958333 | 1.000000 | 0.951456 | 0.970588 |
| 1 | Recall | 0.993589 | 0.931746 | 0.910891 | 0.991150 | 1.000000 | 0.990000 |
| 2 | F1 | 0.994607 | 0.895500 | 0.934010 | 0.995556 | 0.975124 | 0.980198 |
| 3 | Count | 18561 | 630 | 101 | 113 | 196 | 200 |

20 % of dataset left for Validation

**Maastricht University**

# Results and Analysis — Explainability. LIME

- local interpretation
- local, more interpretable model
- approximates underlying model in sample neighborhood
- trained on corrupted copies of original sample
- detects important features in the data

# Results and Analysis - Explainability. LIME

# Results and Analysis - Explainability. LIME

# Results and Analysis - Explainability. LIME

**y=i-Societal Stereotypes** (probability **0.989**, score **4.514**) top features

| Contribution? | Feature |
|---|---|
| 9.708 | Highlighted 9.708 score |
| -5.194 | <BIAS> |

Forward landlord updates regarding events/activities: repairs, holiday/weekend/late night hours, security, parking, shuttle services etc.

**y=i-Societal Stereotypes** top features

| Weight? | Feature |
|---|---|
| +7.688 | [1] landlord |
| +0.623 | [7] holiday |
| +0.620 | [0] Forward |
| +0.418 | [6] repairs |
| +0.366 | [5] activities |
| +0.326 | [14] shuttle |
| +0.152 | [15] services |
| +0.118 | [9] late |
| +0.098 | [10] night |
| +0.058 | [4] events |
| +0.056 | [3] regarding |
| +0.024 | [12] security |
| -0.082 | [8] weekend |
| -0.732 | [11] hours |
| -5.194 | <BIAS> |

# Future Work

# Future work - Suggestion System

1.The **biased words and phrases** that have already been identified in a job description.

⬇

2.**Masked-Language Modeling**: the model will **mask** the biased terms and then give possible **alternatives** according to neighboring words.

⬇

3.Filter out the biased alternatives in the **bias list** or are **classified as bias** in our identification system.

⬇

4.**Order** the unbiased alternatives based on some **criteria** such as cosine similarity.

⬇

5.Finally, the user will choose the most satisfactory alternative from the list.

**Maastricht University**

# Future work - Other Bias Types

- We only focused on five bias types:
  - Generic He/Generic She
  - Societal Stereotypes
  - Behavioural Stereotypes
  - Explicit Marking of Sex
- Other bias types from Jad Doughman's Gender Bias Taxonomy like Benevolent Sexism are also possible to appear in the job descriptions.
- In addition, it is necessary to read more related papers to expand the bias types.

| Bias Type | Bias Subtype | Example | Methodology |
|---|---|---|---|
| Generic Pronouns | Generic He | A programmer must carry his laptop with him to work. | Supervised Learning |
| | Generic She | A nurse should ensure that she gets adequate rest. | Supervised Learning |
| Stereotyping Bias | Societal Stereotypes | Senators need their wives to support them throughout their campaign. | Supervised Learning |
| | Behavioural Stereotypes | The event was kid-friendly for all the mothers working in the company. | Supervised Learning |
| Sexism | Hostile Sexism | Women are incompetent at work. | Supervised Learning |
| | Benevolent Sexism | They're probably surprised at how smart you are, for a girl. | Supervised Learning |
| Exclusionary Terms | Explicit Marking of Sex | Chairman, Businessman, Manpower, Cameraman | Lexicon-Based |
| | Gendered Neologisms | Man-bread, Man-sip, Man-tini | Lexicon-Based |

Figure 2: Overview of the taxonomy and link to detection methodology ([2])

**Maastricht University**

# Future work - Modeling

"It is necessary a dedicated and **ambitious** food development director to join our team."

dedicated and **ambitious** food development

- Experiment with different context size
- Use padding if the biased term is in the beginning or the end of the sentence
- Use weights for the context vs the biased term

**Maastricht University**

# References

[1]  R.Frissen.A-machine-learning-approach-to-recognize-bias-and-discrimination-in-job-advertisements, 2021

[2]  J.Doughman and W.Khreich.Gender bias in text:Labeled datasets and lexicons.CoRR, abs/2201.08675, 2022

[3] Employment Scam Aegean Dataset, 01 2020

[4] Text Analytics Explained-Job Description Data, 09 2018.https://www.kaggle.com/datasets/airiddha/trainrev1

**Maastricht University**

# Thank you for your attention