SOS: Systematic Offensive Stereotyping Bias in Word Embeddings

Anonymous ACL submission

Abstract

Hate speech detection models aim to provide a safe environment for marginalised social groups to express themselves. However, the bias in these models could lead to silencing those groups. In this paper, we introduce the systematic offensive stereotyping (SOS) bias metric. We propose a method to measure the SOS bias in different word embeddings and also investigate its influence on the downstream task of hate speech detection. Our results show that SOS bias against various groups exists in widely used word embeddings and that, in most cases, our SOS bias metric 014 correlates positively with the bias statistics of published surveys on online abuse and hate. 016 However, we found that it is not easy to prove that bias in word embeddings influences downstream task performance. Finally, we show that our SOS bias metric is more indicative of sexism and racism in the inspected word embeddings when used for sexism and racism detection than the stereotypical social biases.

1 Introduction

017

034

040

Wagner et al. (2021) describe the term algorithmically infused societies as the societies that are shaped by algorithmic and human processes and behaviour. The data that is collected from these algorithmically infused societies carry the same bias in algorithms and humans, like population bias and behavioural bias (Olteanu et al., 2019). Among the algorithmically infused societies are social media platforms like Twitter (Dorsey et al., 2021), Urban Dictionary (Peckham, 2021) and 4chan (Poole, 2021); collaborative platforms like Wikipedia (wales and Sanger, 2021) and news aggregating platforms like Google-news (google, 2021). These platforms have various biases. For example, social media platforms have been shown to be rife with offensive and racially insensitive comments (Nguyen et al., 2017; Voué et al., 2020; Mittos et al., 2020). For Wikipedia, in addition

to having language biases (Miz et al., 2020), individuals' biases can be translated into a collective 043 bias (Oeberst et al., 2016), and the news covered in 044 Google News has been shown to be skewed toward 045 the US and the EU in both English and non-English 046 news (Segev, 2008). These biases are important 047 in the field of Natural Language Processing (NLP) because unsupervised models like word embeddings encode them during training. (Brunet et al., 2019; Joseph and Morgan, 2020). This includes 051 racial biases (Garg et al., 2018; Manzini et al., 052 2019; Sweeney and Najafian, 2019), gender biases (Garg et al., 2018; Bolukbasi et al., 2016; Chaloner and Maldonado, 2019), and personality stereotypes (Agarwal et al., 2019). However, one aspect of 056 bias that has received less attention is systematic 057 offensive stereotyping in word embeddings, which includes associating offensive terms to different groups of people, especially marginalised people, 060 based on their ethnicity, gender, or sexual orienta-061 tion. On the other hand, studies that focused on the 062 same bias in hate speech detection models studied 063 it within hate speech datasets (Dixon et al., 2018; 064 Waseem and Hovy, 2016a; Zhou et al., 2021), but 065 not in the widely-used word embeddings which are, 066 in contrast, not trained on data specifically curated 067 to contain offensive content. Moreover, most of 068 the studies on bias in word embeddings focused on 069 studying bias in Word2Vec (Mikolov et al., 2013) 070 and GloVe (Pennington et al., 2014). However, re-071 cent pre-trained word embeddings models like the 072 Urban Dictionary word embeddings that were pre-073 trained on words and definitions from the Urban 074 Dictionary website (Wilson et al., 2020), the Chan 075 word embeddings that were pre-trained on 4& 8 Chan websites (Voué et al., 2020), and a version 077 of GloVe pre-trained on Twitter data (Stojanovski 078 et al., 2015) have received much less attention in 079 previous studies of bias. As we have previously noted, the social media platforms on which these 081 embeddings have been trained are biased (Nguyen

et al., 2017; Voué et al., 2020; Mittos et al., 2020; Mislove et al., 2011). Additionally, the literature on bias in word embeddings claims that it influences downstream tasks, like translation, text classification, and text generation. Still, these claims have not yet been tested (Blodgett et al., 2020).

084

091

097

100

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

122

123

124

125

126

127

129

130

131

132

133

In this work, we are interested in answering the following research questions: **RQ1**: Do word embeddings have systematic offensive stereotyping (SOS) bias, and how can we measure it? **RQ2**: Among the examined word embedding models, which has the most SOS bias? How reflective is SOS bias to online hate and abuse? **RQ3**: How does SOS bias in word embeddings relate to performance on downstream tasks? **RQ4**: How does SOS bias differ from stereotypical social bias regarding finding the most biased word embeddings when used for the task of hate speech detection?

To answer RQ1, we built on the existing literature on measuring bias in word embeddings and proposed a method to measure SOS bias in word embeddings by investigating how different word embedding models associate profanity with marginalised groups of people. To answer the first part of RQ2, we computed the SOS bias score for five different word embedding models and compared their scores, and to answer the second part of the RQ2, we compared our SOS bias scores of the different word embeddings to online surveys on online abuse. To answer RQ3 and to understand how the SOS bias in word embeddings influences downstream task performance, we consider the following tasks: (a) offensive words categorisation and (b) hate speech detection, and measure how performance on these tasks correlates with SOS bias. Finally, to answer RQ4 and find out whether SOS or stereotypical social bias, as measured by stateof-the-art metrics, is more indicative of the bias in the examined word embeddings for the task of hate speech detection, we investigated which bias metrics correlate with the F1 scores of deep learning models using the different word embeddings trained and tested on hate speech related datasets.

The contributions of this paper can be summarised as follows: (a) We define the SOS bias, propose a method to measure it in word embeddings, and demonstrate that our SOS metric results are, for some word embeddings, representative of the abuse that marginalized people experience online and in line with published statistics on online abuse. (b) We demonstrate that all the examined word embeddings contain SOS bias, regardless of the source of the data that they were trained on, with variations on the strength of the bias towards one particular marginalised group or another. (c) We demonstrate that the claim that bias in word embeddings influences downstream tasks is not easy to prove and that despite finding a positive correlation between the SOS bias results and the performance on the downstream tasks, it is not conclusive. (d) We demonstrate that the SOS metric is more indicative of the sexism and racism in the inspected word embeddings than the stereotypical social bias, gender, and racial biases, as measured by state-of-the-art metrics when used for the task of hate speech detection. 134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

183

Our findings show that the different word embeddings contain offensive bias, particularly towards marginalised groups, and it does have an influence, to some extent, on the downstream tasks of hate speech and abuse detection. This bias could have negative implications as these hate speech detection models might learn to associate marginalised groups with hate and abuse. As a result, these models that were supposed to provide a protective environment for the marginalised people to express themselves are the ones that could lead to silencing them or flagging their content as inappropriate.

2 Background: Bias and Word Embeddings

The term bias is defined and used in many different ways (Olteanu et al., 2019). There is the normative definition of bias, as its definition in cognitive science as: "behaving according to some cognitive priors and presumed realities that might not be true at all" (Garrido-Muñoz et al., 2021). There is also the statistical definition of bias as "systematic distortion in the sampled data that compromises its representatives" (Olteanu et al., 2019). In the literature on bias in word embeddings, we find different definitions for bias. For example, Caliskan et al. (2017) define bias from a normative perspective as prior information necessary to requisite for intelligent action, while (Sweeney and Najafian, 2019) define bias from a statistical perspective as unequal distribution of negative sentiment among demographic identity terms in word embeddings. Similarly, (Dev and Phillips, 2019) defines bias as deviation from a population parameter.

In the case of distributional word representations (Word Embeddings), the literature had focused on

commonly used word embeddings (Elsafoury et al., 2021a), like Word2Vec (Mikolov et al., 2013) and 185 GloVe (Pennington et al., 2014). However, more recent word embeddings like urban dictionary (UD) (Wilson et al., 2020) and Chan (Voué et al., 2020) have not been well studied for bias, even though there is evidence from the literature that the data that was used in pre-training these word embeddings contain offensiveness and racial comments 192 (Nguyen et al., 2017; Voué et al., 2020; Mittos et al., 2020; Mislove et al., 2011).

184

187

189

190

191

193

195

196

197

199

201

204

210

211

212

213

214

215

216

217

219

227

The most common methods for quantifying bias in word embeddings are WEAT, RND, RNSB, and ECT. For WEAT, the authors were inspired by the Implicit Association Test (IAT) to develop a statistical test to demonstrate human-like biases in word embeddings (Caliskan et al., 2017). They used the cosine similarity and statistical significance tests to measure the unfair correlations for two different demographics, as represented by manually curated word lists. For RND, the authors used the Euclidean distance between neutral words, like professions, and a representative group vector created by averaging the word vectors for words that describe a stereotyped group (gender/ethnicity) (Garg et al., 2018). In RNSB, a logistic regression model was first trained on the word vectors of unbiased labeled sentiment words (positive and negative) extracted from biased word embeddings. Then, that model was used to predict the sentiment of words that describe certain demographics (Sweeney and Najafian, 2019). In ECT, the authors proposed a method to measure how much bias has been removed from the word embeddings after debiasing them (Dev and Phillips, 2019).

These measures, except RNSB, are based on the polarity between two opposing points, like male and female, allowing for binary comparisons. This forces practitioners to model gender as a spectrum between more "male" and "female" words, requiring an overly simplified view of the construct, leading to similar problems for other stereotypical types of bias, like racial, religious, transgender, and sexual orientation, where there are more than two categories that need to be represented (Sweeney and Najafian, 2019). These metrics also use lists of seed words that have been shown to be unreliable (Antoniak and Mimno, 2021). Since we are interested in measuring the systematic offensive stereotypes of different marginalised groups, this measure would fall short of our needs. As for the RNSB measure, even though it is possible to include more than two identities, the sentiment dimension is represented as positive or negative (binary). But in our case, we are interested in a variety of offensive language targeted at different marginalised groups.

Systematic Offensive Stereotyping Bias 3

Our motivation is to reveal whether word embeddings associate offensive language with words describing marginalised groups. We define systematic offensive stereotyping (SOS) bias from a statistical perspective as "A systematic association in the word embeddings between profanity and marginalised groups of people". In the next section, we will use this definition to measure the SOS bias and to answer RQ1.

3.1 Measuring SOS bias

In this section, we describe our proposed method to measure SOS bias in various word embeddings. Based on our definition of SOS, we want a method to measure the association that each word embedding model has between profanity and marginalised groups of people. We propose to measure that association using the cosine similarity between swear words and words that describe marginalised social groups. For the swear words, we use a list of 427

Group	Word				
LGTBQ*	lesbian, gay, queer, homosexual, lgbt, bi-				
	sexual, transgender, trans, non-binary				
Women*	woman, female, girl, wife, sister, mother,				
	daughter				
Other ethnicities*	african, african american, black, asian, his-				
	panic, latin, mexican, indian, arab				
Straight	hetrosexual, cisgender				
Men	man, male, boy, son, father, husband,				
	brother				
White ethnicities	white, caucasian, european american, eu-				
	ropean, norwegian, canadian, german, aus-				
	tralian, english, french, american, swedish,				
	dutch				

*Marginalised group

Table 1: NOI words and the group they describe.

swear words collected by (Agrawal and Awekar, 2018). For describing marginalised social groups, we used a word list that contains non-offensive identity (NOI) names to describe marginalised groups of people (Zhou et al., 2021; Dixon et al., 2018) and non-marginalised ones (Table 1).

Let $W_{NOI} = \{w_1, w_2, w_3, ..., w_n\}$ be the list of NOI words w_i , i = 1, 2, ..., n, and $W_{sw} =$ $\{o_1, o_2, o_3, \dots o_m\}$ be the list of swear words o_i , j = 1, 2, ..., m. To measure the SOS bias for a

267

269

259

237

238

239

240

241

242

243

244

245

246

247

248

250

251

252

253

254

255

256

257



Figure 1: Mean SOS scores for the examined word embeddings and groups.

specific word embedding we, we first compute the average vector $\overline{\mathbf{W}_{sw}^{we}}$ of the swear words for we, e.g. for Word2Vec, Glove, etc. $SOS_{i,we}$ for a NOI word w_i and a word embedding we is then defined (Equation 1) as the cosine similarity between $\overline{\mathbf{W}_{sw}^{we}}$ and the word vector $\overline{w_{i,we}}$, for the word embedding we, normalised to the range [0, 1] using min-max normalisation.

270

273

276

277

278

281

284

285

286

292

296

297

301

$$SOS_{i,we} = cos(\overrightarrow{\mathbf{w}_{sw}^{we}}, \overrightarrow{w_{i,we}}) = \frac{\overrightarrow{\mathbf{w}_{sw}^{we}} \cdot \overrightarrow{w_{i,we}}}{||\overrightarrow{\mathbf{w}_{sw}^{we}}|| \cdot ||\overrightarrow{w_{i,we}}||}$$
(1)

The normalised SOS score takes values within the range [0, 1] and indicates the similarity of a NOI word to the average representation of swear words. Consequently, a higher $SOS_{i,we}$ value for word w_i indicates that the word embedding $\overrightarrow{w_{i,we}}$ for the word w_i , is more associated with profanity.

3.2 Mean SOS for word embeddings

We then proceeded to compute the mean SOS score for the following five word embeddings: Word2Vec, Glove-WK, Glove-Twitter, UD, and Chan, using the aforementioned swear words and NOI word lists for each examined group individually, as well as for the combined marginalised (Women, LGBTQ, Other ethnicities) and nonmarginalised (Men, Straight, White ethnicity) groups. Figure 1 answers RQ1, showing that there is SOS bias in the word embeddings towards all the examined groups, both marginalised and nonmarginalised. In addition, Table 2 shows that mean SOS bias towards the marginalised groups is higher than towards the non-marginalised groups (T-test p = 0.02 for $\alpha = 0.05$).

It is also evident that when comparing the "Straight" and the "LGBTQ" groups, there is a higher SOS bias towards the marginalised "LGBTQ" group for all the examined word em-

Word ombodding	Mean SOS			
word embedding	Marginalised	Non-marginalised		
Word2Vec	0.535	0.430		
Glove-WK	0.390	0.281		
Glove-Twitter	0.558	0.469		
UD	0.407	0.325		
Chan	0.495	0.417		

Table 2: Mean SOS score of the different groups.

beddings. Similar for the "Men" vs. "Women" groups and "White ethnicity" vs. "Other ethnicities" groups, where there is higher SOS bias towards the marginalised "Women" and "Other ethnicities" groups, except for Glove-WK and UD for which the SOS bias is marginally higher for the nonmarginalised groups ("Men", "White ethnicity"). The rest of this work will focus on the marginalised groups (women, LGBTQ, other ethnicities).

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

328

330

331

332

333

334

335

336

337

338

339

340

342

343

344

345

346

3.3 Which word embedding has the most SOS bias?

To answer the first part of RQ2, we conducted a comparative analysis between the word embeddings in regards to SOS bias. To quantitatively compare the different word embeddings, we used the SOS bias scores (Figure 1) for each marginalised group (LGTBQ, Women, Other ethnicities) and applied the Friedman and T-test significance tests $(\alpha = 0.05)$. For the words that describe the "LGTBQ" group, Glove-WK has the highest SOS score of 0.629, but the Friedman test failed in finding a significant difference between the different word embeddings (p = 0.6), indicating that all the examined word embeddings are similarly SOS-biased towards words related to the "LGBTQ" group. For the "Women" group, Glove-Twitter, UD, and Chan exhibited high SOS bias, with Glove-Twitter having the highest score of 0.852, and Friedman's test indicating a significant difference between the word embeddings $(p = 5e^{-4})$. A T-test showed that Glove-Twitter is significantly different from Word2Vec, Glove-WK, and UD $(p = 6e^{-6}, 1e-5, and 0.0057 respec$ tively) but no significant difference from Chan (p =0.350) could be established. This indicates that Glove-Twitter and Chan exhibit a similar significant SOS bias towards women (sexism) in comparison to Word2Vec, Glove-WK, and UD. Regarding the "Other ethnicities" group, Word2Vec stands out as the word embedding with the highest SOS score of 0.691. Friedman's test showed a statistically significant difference between all

Word Embedding	SOS biased towards
Word2Vec	Other ethnicities, LGBTQ, Women
Glove-WK	LGBTQ, Women, Other ethnicities
Glove-Twitter	Women, Other ethnicities, LGBTQ
UD	Women, LGBTQ, Other ethnicities
Chan	Women, LGBTQ, Other ethnicities

Table 3: The groups that each word embedding is SOS-biased towards, ordered by descending severity.



Figure 2: The Spearman's rank correlation coefficient between the ranking of SOS measure and the ranking of the mean collocation PMI.

the word embeddings (p = 4e-4) and the T-test showed that the SOS score of Word2Vec is significantly higher than Glove-WK, Glove-Twitter, UD, and Chan $(p = 9e^{-7}, 8e^{-3}, 1e^{-5}, and 4e^{-5})$ respectively), indicating that Word2Vec is significantly SOS-biased towards non-white ethnicities in comparison to Glove-WK, Glove-Twitter, UD, 353 and Chan. We summarise our results in Table 3 showing that Word2Vec is the most SOS-biased towards non-white ethnicities, Glove-WK is the most SOS-biased towards the LGBTQ community, and Glove-Twitter, UD, and Chan are the most SOSbiased towards women.

3.4 **Results Validation**

347

351

354

357

362

364

366

367

370

371

To validate our results, and to answer the second part of RQ2, we compared our results to the collocations between the NOI words of marginalised groups and swear words following the work of (Pietraszewska, 2013). To generate these collocations, we used a corpus of randomly sampled 100,000 Pushshift's public Reddit collection (Reddit, 2021) comments (4 million tokens) that were posted between 2005 and 2012. Then, we used NLTK (NLTK, 2021) to find the words that cooccur the most with the NOI words and filtered them to find the co-occurrences between the NOI words w_i and the swear words o_i . The association between the acquired word pairs was measured

using the pointwise mutual information (PMI). Then we computed the mean PMI for all the cooccurrences of offensive words and each of the NOI words (Equation 2). Finally, we computed the Spearman's rank correlation coefficient between the ranked mean PMI, PMI_i , and the ranked SOS score $SOS_{i,we}$, for each NOI word w_i and word embedding we.

$$\overline{PMI_i} = \frac{1}{m} \sum_{j=1}^{m} PMI(w_i, o_j)$$
(2)

375

376

377

378

379

380

381

383

385

387

388

389

390

391

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

Results in Figure 2, show a positive correlation for all the marginalised groups and most of the word embeddings, except for Glove-WK for "Other ethnicities" and Word2Vec, UD, and Chan for "Women", where a negative correlation is detected. After inspecting the "Women"-related words, where the correlation is negative, we found that they collocated with slurs that are not widely used and were not included in the used swear words list**. The correlation for all the NOI words in the marginalised group shows a positive correlation with all the word embeddings except for Glove-WK. We speculate that this is the case because, as shown in Figure 1 and Table 3, Glove-WK is the least biased towards "Other ethnicities".

In addition to the collocations, we computed the Pearson's correlation coefficient between our SOS metric and metrics from two published surveys of online abuse. The surveys are the Rad Campaign Online Harassment Survey 2014 (Rad Campaign, 2014) where 1,000 adult Americans (aged 18+) were surveyed about being harassed online. The other survey is the COX Teen Internet Safety Survey (Cox Communications Inc., 2014), where a total of 1,301 teens aged 13-17 were surveyed about being bullied online. We chose these two surveys because they provide data on all the marginalised groups that we analyse in this paper (women, other ethnicities, and LGBTQ). Results in Table 4 show a positive correlation between the SOS metric and the surveys' metrics for Glove-Twitter, UD, and Chan. However, there is also a strong negative correlation between Word2Vec, Glove-WK, and the surveys' metrics. We also used the survey data on online extremism and online hate (OEOH), collected by (Hawdon et al., 2015) from Finland (n=555), the US (n=1,033), Germany

^{**}We have not added these slurs to the swear words' list as more validation work would be required to confirm that they unambiguously belong in the list, thus risking biasing our results based on our own observations.

SUPPON	Pearson's correlation							
Survey	Word2Vec	Glove-WK	Glove-Twitter	UD	Chan			
Rad*	-0.608	-0.486	0.977	0.344	0.666			
COX	-0.964	0.142	0.905	0.843	0.981			
OEOH-US	0.846	0.151	-0.989	-0.650	-0.884			
OEOH-Germany	0.842	0.158	-0.990	-0.644	-0.881			
OEOH-Finland	0.912	0.013	-0.960	-0.749	-0.940			
OEOH-UK	0.934	-0.043	-0.943	-0.785	-0.958			
*Completion computed between SOS and the differences in Ded between the nerventees								

of (women and men), (Other ethinicities and Caucasian), and (LGBTQ and straight).

Table 4: correlation between survey metrics and SOS.

(n=978), and the UK (n=999) in 2013 and 2014. The respondents are individuals aged (15 - 30). The correlation shows almost the opposite pattern to the correlation with surveys on online abuse, as Word2Vec and Glove-WK correlate positively with online extremism and hate, while Glove-twitter, UD, and Chan have a strong negative correlation to online extremism and hate for the US and Germany. For the UK and Finland, Word2Vec shows a stronger positive correlation with the online extremism data, but Glove-WK shows no correlation.

421

422

423

424

425

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

These results suggest that the word embeddings that were trained on the social media datasets (Glove-twitter, UD, and Chan) encode the online abuse towards marginalised people, while word embeddings that were trained on Google news and Wikipedia articles encode the hate and extremism against the marginalised groups shared in those sources. These results reveal that our SOS metric correlates with other measures of hate and abuse towards the marginalised groups in social media data. However, results are not conclusive, and more in-depth analysis with more datasets is required.

4 SOS bias and downstream tasks

In this section, we answer RQ3 through a series of experiments on one downstream task, i.e., hate speech detection. A second downstream task, offensive words categorisation, is also examined in Appendix A.1. We investigated the influence of SOS bias in the word embeddings on the task of hate speech detection by training deep learning models with an embedding layer for the detection of different types of hate speech from hate speechrelated datasets, then computed the correlation of the performance of the different word embeddings to the SOS bias score of these embeddings.

4.1 Datasets & Pre-processing

We used four hate speech-related datasets from social media sources that contain different types

Dotocot	Somplos	Positive	Avg. words	Max. words		
Dataset	Samples	samples	per comment	per comment		
HateEval	12722	42%	21.75	93		
Twitter-sexism	14742	23%	15.04	41		
Twitter-racism	13349	15%	15.05	41		
Twitter-hate	5569	25%	14.60	32		
Note: Positive samples refer to offensive comments						

Table 5: Hate speech datasets' statistics.

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

of hate speech: (i) *Twitter-racism*, a collection of Twitter messages containing tweets that are labeled as racist or not (Waseem and Hovy, 2016b); (ii) *Twitter-sexism*, Twitter messages containing tweets labeled as sexist or not (Waseem and Hovy, 2016b); (iii) Twitter-hate, containing tweets that are labeled as offensive, hateful (sexist, homophobic, and racist), or neither (Davidson et al., 2017). As we are interested in the hateful content, we used the tweets that are labeled as hateful or neither; and (iv) *HateEval*, a collection of tweets containing hate speech against immigrants and women in Spanish and English (Basile et al., 2019), from which we used only the English tweets. Statistics about the datasets are provided in Table 5.

To pre-process the datasets, we removed URLs, user mentions, retweet abbreviation "RT", non-ASCII characters, and English stop words except for second-person pronouns like "you/yours/your", and third-person pronouns like "he/she/they", "his/her/their" and "him/her/them" were not removed, as suggested in (Elsafoury et al., 2021b). All letters were lowercased, and common contractions were converted to their full forms. Finally, each dataset was randomly split into training (70%) and test (30%) sets, preserving class ratios.

4.2 Machine Learning Models

We used two deep learning models: (i) a Bidirectional LSTM (Schuster and Paliwal, 1997) with the same architecture as in (Agrawal and Awekar, 2018), who used RNN models to detect hate speech, and (ii) a two layers Multi-Layer Perceptron (MLP) model. To this end, we first used the Keras tokenizer (Tensorflow.org, 2020) to tokenise the input texts, using a maximum input length of 64 (maximum observed sequence length in the dataset). A frozen embedding layer, based on a given pretrained word embedding model, was used as the first layer and fed to the BiLSTM model and to the MLP model. To avoid over-fitting, we used L2 regularisation with an experimentally determined value of 10^{-7} . For each dataset, the models were 502 503

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

524

526

527

528

trained on the training set for 100 epochs with a batch size of 32, using the Adam optimiser and a learning rate of 0.01 (default of Keras Optimiser).

Datacat	Model	F1-score						
Dataset	wouer	Word2Vec	Glove-WK	Glove-Twitter	UD	Chan		
HotoEvol	MLP	0.593	0.583	0.623	0.597	0.627		
HateLvai	BiLSTM	0.663	0.651	0.671	0.661	0.661		
Twitter-sexism	MLP	0.587	0.587	0.589	0.578	0.563		
	BiLSTM	0.659	0.661	0.661	0.625	0.631		
Twitter region	MLP	0.683	0.681	0.680	0.679	0.650		
1 witter-racisin	BiLSTM	0.717	0.727	0.6999	0.698	0.712		
Twitter-hate	MLP	0.681	0.713	0.775	0.780	0.692		
	BiLSTM	0.772	0.821	0.851	0.837	0.84		

Note: Numbers in bold indicate best performance per model and dataset

Table 6: F1 scores for the used models using the examined word embeddings on our datasets.

Dataset	Model	ρ
H-4-E1	MLP	0.500
HateEval	BiLSTM	0.974
Twitter-sexism	MLP	0.461
	BiLSTM	0.205
Twitter-racism	MLP	0.1
	BiLSTM	-0.3
Twitter-hate	MLP	-0.2
	BiLSTM	0.4

Table 7: Spearman's correlation between mean SOS and the F1 scores for the used models using the examined word embeddings for our datasets.

4.3 Experimental Results

Given the results for the SOS bias in the different embeddings (Table 3), we hypothesise that the deep learning models that are trained with Word2Vec embeddings will perform the best (highest F1 score) on datasets that contain hate speech or insults towards marginalised ethnicities, which is Twitterracism. We also hypothesise that the models trained with Glove-Twitter, UD, and Chan will achieve the highest F1 scores on datasets that contain insults towards women, which are Twitter-racism and HateEval. Given that the Twitter-hate dataset contains a mixture of sexist, homophobic, and racist comments, we hypothesise that the models trained with Glove-Twitter, UD, and Chan will perform the best.

The classification performance of the deep learning models with the different embedding models is reported in Table 6. The results show that for all datasets, BiLSTM outperforms MLP in terms of F1 score. In addition, results show that for the MLP model, our hypotheses hold for all four datasets, as Chan is the best performing for a dataset that contains insults towards women (HateEval), Word2Vec is the best performing on a dataset that contains insults towards other ethnicities (Twitter-racism), Glove-Twitter is the best performing on a dataset that contain insults towards women (Twitter-sexism), and UD is the best performing on Twitter-hate which contain insults towards women and the LGBTQ community. For the BiLSTM model, our hypotheses hold for three datasets, i.e., HateEval, Twitter-sexism, and Twitter-hate, as Glove-Twitter is the best performing on datasets that contain insults towards women and LGTBQ, which are found in the HateEval, Twitter-sexism, and Twitter-hate datasets. As for the Twitter-racism dataset, we hypothesised that Word2Vec would be the best performing, but instead, Glove-WK is the best performing when the BiLSTM model is used. 531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

To quantify our analysis of the influence of the SOS bias on the task of hate speech detection, we used Spearman's rank correlation coefficient to compute the correlation between the ranking of the mean SOS bias scores and the ranking of F1 scores for the MLP and BiLSTM models for the different word embeddings in each examined datasets. As shown in Table 7, there is a positive correlation between the mean SOS bias score and the BiLSTM F1 scores of the different word embeddings for Twitter-sexism, HateEval, and Twitter-hate, but not for the Twitter-racism dataset where the correlation is negative. For the MLP models, the correlation was positive for three datasets, i.e., Twitter-sexism, Twitter-racism, and HateEval, but not for Twitterhate. The results in this section and in Appendix A.1 suggest that SOS bias in word embeddings influences the performance of downstream tasks. This finding is less evident for the task of offenses categorisation but clearer for the task of hate speech detection. However, results are not conclusive and more experiments are required.

5 Comparative analysis of bias metrics

To answer RQ4, we compared our SOS bias metric to state-of-the-art bias metrics from the literature, in particular WEAT, RND, RNSB, and ECT, regarding finding the most biased word embeddings for the task of hate speech detection. In this section, we performed the comparison on the task of sexism detection. Thus the metrics were used to measure gender bias. The same experiment was also conducted for racial bias in Appendix A.2. We used the WEFE framework (Badilla et al., 2020) to measure the bias using the other state-of-the-art metrics.

To measure the gender bias in the word em-

beddings using the state-of-the-art metrics, we 581 used two target lists: Target list 1, which con-582 tains female-related words (e.g., she, woman, and 583 mother), and Target list 2, which contains malerelated words (e.g., he, father, and son), and two attribute lists: Attribute list 1, which contains words 586 related to career, science, math, intelligence, male 587 roles, as well as positive words, and Attribute list 2, which contains words related to family, arts, appearance, sensitivity, female roles, as well as negative words (Badilla et al., 2020; Caliskan et al., 2017). 591 Then, we measured the average gender bias scores 592 across the different attribute lists for each word em-593 bedding using the different state-of-the-art metrics. 594 For the SOS bias, we used the mean SOS scores of the words that belong to the "Women" category, as computed in Section 3.2 (Figure 1).

> For each bias metric, we ranked the bias scores for each word embedding in ascending order, except for the ECT metric that was ranked in descending order, as the higher the value, the lower the bias. We then computed the Spearman's rank correlation coefficient between the gender bias of the different word embeddings as measured by WEAT, RND, RNSB, ECT, SOS_{women}), and the F1 scores achieved by the two deep learning models on the Twitter-sexism, HateEval, and Twitter-hate datasets using the different word embeddings (as computed in Section 4.3/Table 6). The computed Spearman's correlations are shown in Table 8.

601

603

607

609

610

611

612

613

614

615

616

617

618

621

625

627

628

629

631

Our results show that for HateEval and Twitterhate, SOS_{women} has a higher positive correlation to the F1 scores of the deep learning models than the rest of the bias metrics, indicating that the SOS bias score of the different word embeddings correlates positively with the performance of the deep learning models using the word embeddings for the task of hate speech detection on these two datasets. However, for Twitter-sexism, SOS_{women} shows almost no correlation with the F1 scores of either MLP or BiLSTM. We speculate that the reason is that 66% of the Twitter-sexism dataset contains sexist tweets that are not profane, in comparison to only 40% in HateEval and Twitter-hate datasets.

Our analysis showed that the gender bias scores of WEAT, ECT, RND, and RNSB metrics for the different word embeddings do not always correlate with the deep learning models' performances using the same word embeddings on the genderrelevant datasets and differs drastically from one dataset to another. The proposed SOS bias score for the different word embeddings shows a more consistent positive correlation with the F1 scores of the deep learning models using these word embeddings when profanity is used against the biastarget group. Similar results were found for racial bias, as presented in Appendix A.2. This indicates that our proposed SOS bias metric is more inductive of the sexist and racist word embeddings than the stereotypical social bias, as measured by the state-of-the-art metrics when used for hate speech detection. 632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

Dataset	Model	Spearman's correlation				
	WIGHT	WEAT	RNSB	RND	ECT	SOS
HateEval	MLP	-0.600	0.300	0.300	0.600	0.800
	BiLSTM	-0.410	-0.718	-0.307	0.666	0.359
Twitter-sexism	MLP	0.153	-0.102	-0.205	0.35	0.051
	BiLSTM	0.564	0.461	0.359	0.416	0.05
Twitter-hate	MLP	-0.700	0.100	-0.400	-0.300	0.500
	BiLSTM	-0.600	0.300	0.300	0.600	1

Table 8: Spearman's rank correlation coefficient of the gender bias scores of the different word embeddings and the F1 scores of the used models for each bias metric and dataset.

6 Conclusion

In this work, we introduced the SOS bias and proposed methods to measure it, validate it, investigate its influence on downstream tasks, and compare it to stereotypical social bias. Our results show that there is SOS bias in the examined word embeddings and that for some of them, it has a strong positive correlation with published statistics on online abuse and extremism. However, more datasets need to be collected to provide stronger evidence, especially data from social sciences on the offenses that marginalised groups receive on social media. Our findings show that proving the influence of bias in word embeddings on the downstream tasks is not an easy task and that even though our results suggest that there is a relationship between the SOS bias and the downstream task of hate speech detection, the results are not conclusive, as there might be other factors that contributed to the performance of the examined deep learning models. As future work, more experiments are required using counterfactual datasets and feature importance scores of NOI words to ensure that we understand the impact of the SOS bias in the word embeddings on the downstream tasks. Finally, our findings suggest that our proposed SOS bias metric is more indicative of the biased word embeddings in comparison to stereotypical social bias for the tasks of sexism and racism detection.

References

672

673

676

679

681

700

701

702

703

705

706

707

708

710

712

714

715

716

717

718

719

720

721

722

723

724

727

- Oshin Agarwal, Funda Durupinar, Norman I. Badler, and Ani Nenkova. 2019. Word embeddings (also) encode human personality stereotypes. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics (*SEM 2019)*, pages 205– 211, Minneapolis, Minnesota. Association for Computational Linguistics.
 - Sweta Agrawal and Amit Awekar. 2018. Deep learning for detecting cyberbullying across multiple social media platforms. In Advances in Information Retrieval - 40th European Conference on IR Research, ECIR 2018, Grenoble, France, March 26-29, 2018, Proceedings, volume 10772 of Lecture Notes in Computer Science, pages 141–153. Springer.
 - Maria Antoniak and David Mimno. 2021. Bad seeds: Evaluating lexical methods for bias measurement. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1889–1904, Online. Association for Computational Linguistics.
 - Pablo Badilla, Felipe Bravo-Marquez, and Jorge Pérez. 2020. WEFE: the word embeddings fairness evaluation framework. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 430–436. ijcai.org.
 - Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 54–63, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
 - Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of "bias" in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454– 5476, Online. Association for Computational Linguistics.
 - Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16, page 4356–4364, Red Hook, NY, USA. Curran Associates Inc.
 - Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard S. Zemel. 2019. Understanding the origins of bias in word embeddings. In Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019,

Long Beach, California, USA, volume 97 of Proceedings of Machine Learning Research, pages 803– 811. PMLR. 729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

782

- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Kaytlin Chaloner and Alfredo Maldonado. 2019. Measuring gender bias in word embeddings across domains and discovering new gender bias word categories. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 25–32, Florence, Italy. Association for Computational Linguistics.
- Cox Communications Inc. 2014. 2014 teen internet safety survey. [Online] Accessed 13/9/2021.
- Thomas Davidson, Dana Warmsley, Michael W. Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In Proceedings of the Eleventh International Conference on Web and Social Media, ICWSM 2017, Montréal, Québec, Canada, May 15-18, 2017, pages 512–515. AAAI Press.
- Sunipa Dev and Jeff M. Phillips. 2019. Attenuating bias in word vectors. In The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan, volume 89 of Proceedings of Machine Learning Research, pages 879–887. PMLR.
- Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '18, page 67–73, New York, NY, USA. Association for Computing Machinery.
- Jack Dorsey, Noah Glass, Biz Stone, and Evan Williams. 2021. Twitter. [Online] Accessed 15/9/2021.
- Fatma Elsafoury, Stamos Katsigiannis, Zeeshan Pervez, and Naeem Ramzan. 2021a. When the timeline meets the pipeline: A survey on automated cyberbullying detection. *IEEE Access*, 9:103541–103563.
- Fatma Elsafoury, Stamos Katsigiannis, Steven R. Wilson, and Naeem Ramzan. 2021b. Does BERT pay attention to cyberbullying? In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 1900–1904, New York, NY, USA. Association for Computing Machinery.
- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.

Ismael Garrido-Muñoz, Arturo Montejo-Ráez, Fer-

Google news.

James Hawdon, Atte Oksanen, and Pekka Räsänen.

Kenneth Joseph and Jonathan Morgan. 2020. When do

word embeddings accurately reflect surveys on our

beliefs about people? In Proceedings of the 58th An-

nual Meeting of the Association for Computational Linguistics, pages 4392-4415, Online. Association

Thomas Manzini, Lim Yao Chong, Alan W Black,

and Yulia Tsvetkov. 2019. Black is to criminal

as caucasian is to police: Detecting and removing

multiclass bias in word embeddings. In Proceed-

ings of the 2019 Conference of the North American

Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1

(Long and Short Papers), pages 615-621, Minneapo-

lis, Minnesota. Association for Computational Lin-

Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig.

2013. Linguistic regularities in continuous space

word representations. In Proceedings of the 2013

Conference of the North American Chapter of the

Association for Computational Linguistics: Human

Language Technologies, pages 746-751, Atlanta,

Georgia. Association for Computational Linguistics.

Pekka Onnela, and J. Niels Rosenquist. 2011. Un-

derstanding the demographics of twitter users. In

Proceedings of the Fifth International Conference

on Weblogs and Social Media, Barcelona, Catalo-

nia, Spain, July 17-21, 2011. The AAAI Press.

Alexandros Mittos, Savvas Zannettou, Jeremy Black-

burn, and Emiliano De Cristofaro. 2020. "and we

will fight for our race!" A measurement study of ge-

netic testing conversations on reddit and 4chan. In

Proceedings of the Fourteenth International AAAI

Conference on Web and Social Media, ICWSM 2020,

Held Virtually, Original Venue: Atlanta, Georgia,

USA, June 8-11, 2020, pages 452-463. AAAI Press.

jamin Ricaud, and Pierre Vandergheynst. 2020.

What is trending on wikipedia? capturing trends and

language biases across wikipedia editions. In Com-

panion of The 2020 Web Conference 2020, Taipei,

Taiwan, April 20-24, 2020, pages 794-801. ACM.

Dong Nguyen, Barbara McGillivray, and Taha Yasseri.

2017. Emo, love, and god: Making sense of ur-

Volodymyr Miz, Joëlle Hanna, Nicolas Aspert, Ben-

Alan Mislove, Sune Lehmann, Yong-Yeol Ahn, Jukka-

2015. Online extremism and online hate. NORDI-

Sciences, 11(7).

COM, page 29.

for Computational Linguistics.

google. 2021.

guistics.

15/9/2021.

nando Martínez-Santiago, and L. Alfonso Ureña-

López. 2021. A survey on bias in deep nlp. Applied

[Online] Accessed

- 790 791
- 793
- 796 797

- 802

- 809
- 810 811
- 812 813
- 814
- 815 816
- 817 818

819 820 821

- 822 823 824
- 825 826
- 827
- 829
- 830 831

832 834

838

ban dictionary, a crowd-sourced online dictionary. *CoRR*, abs/1712.08647.

NLTK. 2021. Nltk collocations. https://www. nltk.org/howto/collocations.html. Accessed: 2021-09-14.

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

891

- Aileen Oeberst, Ulrike Cress, Mitja Back, and Steffen Nestler. 2016. Individual Versus Collaborative Information Processing: The Case of Biases in Wikipedia, pages 165–185. Springer International Publishing, Cham.
- Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. Frontiers in Big Data, 2:13.
- Aaron Peckham. 2021. Urban dictionary. [Online] Accessed 15/9/2021.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, pages 1532–1543. ACL.
- Natalia Pietraszewska. 2013. A qualitative and quantitative analysis of selected ethnic and racial terminology present in assorted public english corpora. Styles of Communication, 5(1).
- Christopher Poole. 2021. 4 chan. [Online] Accessed 15/9/2021.
- Rad Campaign. 2014. The rise of online harassment. [Online] Accessed 13/9/2021.
- Reddit. 2021. Pushshift-reddit. https://files. pushshift.io/reddit/comments/. Accessed: 2021-09-14.
- Mike Schuster and Kuldip K Paliwal. 1997. Bidirectional recurrent neural networks. IEEE Transactions on Signal Processing, 45(11):2673-2681.
- Elad Segev. 2008. The imagined international community: Dominant american priorities and agendas in google news. Global Media Journal.
- Dario Stojanovski, Gjorgji Strezoski, Gjorgji Madjarov, and Ivica Dimitrovski. 2015. Twitter sentiment analysis using deep convolutional neural network. In Hybrid Artificial Intelligent Systems, pages 726–737, Cham. Springer International Publishing.
- Chris Sweeney and Maryam Najafian. 2019. A transparent framework for evaluating unintended demographic bias in word embeddings. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1662–1667, Florence, Italy. Association for Computational Linguistics.
- Tensorflow.org. 2020. Text tokenization utility class. https://www.tensorflow.org/api_docs/ python/tf/keras/preprocessing/text/ Tokenizer. Accessed: 2020-09-28.
- 10

- Pierre Voué, Tom De Smedt, and Guy De Pauw. 2020. 4chan & 8chan embeddings. *CoRR*, abs/2005.06946.
- Claudia Wagner, Markus Strohmaier, Alexandra Olteanu, Emre Kıcıman, Noshir Contractor, and Tina Eliassi-Rad. 2021. Measuring algorithmically infused societies. *Nature*, 595(7866):197–204.

896

900

901

902

903

904

905

906

907

908

909

910 911

912

913

914

915

916

917

918

919

920

921

922

924

925 926

927

928

929

930

931

932

933

934

935

938

939

941

943

944

945

- Jimmy wales and Larry Sanger. 2021. Wikipedia. [Online] Accessed 15/9/2021.
- Zeerak Waseem and Dirk Hovy. 2016a. Hateful symbols or hateful people? predictive features for hate speech detection on twitter. In Proceedings of the Student Research Workshop, SRW@HLT-NAACL 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016, pages 88–93. The Association for Computational Linguistics.
 - Zeerak Waseem and Dirk Hovy. 2016b. Hateful symbols or hateful people? predictive features for hate speech detection on Twitter. In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, San Diego, California. Association for Computational Linguistics.
 - Steven R. Wilson, Walid Magdy, Barbara McGillivray, Kiran Garimella, and Gareth Tyson. 2020. Urban dictionary embeddings for slang NLP applications. In Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020, pages 4764–4773. European Language Resources Association.
 - Haoran Zhang, Amy X. Lu, Mohamed Abdalla, Matthew B. A. McDermott, and Marzyeh Ghassemi. 2020. Hurtful words: quantifying biases in clinical contextual word embeddings. In ACM CHIL '20: ACM Conference on Health, Inference, and Learning, Toronto, Ontario, Canada, April 2-4, 2020 [delayed], pages 110–120. ACM.
 - Xuhui Zhou, Maarten Sap, Swabha Swayamdipta, Yejin Choi, and Noah A. Smith. 2021. Challenges in automated debiasing for toxic language detection. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021, pages 3143–3155. Association for Computational Linguistics.

A Appendix

A.1 Offensive words categorisation

We investigated the influence that the SOS bias in the word embeddings has over the downstream task of offenses categorisation. We used the Hurtlex lexicon (Zhang et al., 2020), which is a multilingual lexicon containing 8,228 offensive words and expressions, organised into 17 groups. We used words from the English lexicon that belong to the 11 groups that are related to the marginalised groups studied in this work. The used categories are ethnic slurs (PS); words related to social and economic disadvantage (IS), descriptive words with potential negative connotations (QAS), derogatory words (CDS), felonies and words related to crime and immoral behavior (RE), male genitalia (ASM), female genitalia (ASF), words related to prostitution (PR), words related to homosexuality (OM), cognitive disabilities and diversity (DDP), and physical disabilities and diversity (DDF).

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

To investigate the influence that the SOS bias has on the ability of each word embedding to group together the words that belong to the same Hurtlex category, we trained a KNN model. We first removed the words in the lexicon that belong to more than one category, resulting in 5,963 offensive words in total. We then split the Hurtlex lexicon into a training (70%) and a test (30%) set, preserving the class ratio. The F1-scores achieved by the KNN model for each of the 11 classes for the test set are shown in Figure 3. A Friedman test ($\alpha = 0.05$) between the F1 scores of each data item in the test set showed that the F1 scores achieved using the examined word embeddings are significantly different. To further investigate the difference between pairs of top-scoring word embeddings, we used a Wilcoxon test ($\alpha = 0.05$). Results showed that, across all classes, UD scores significantly higher than Chan and Glove-WK, but not significantly higher than Word2Vec or Glove-Twitter. Similarly, we found that Word2Vec achieves a significantly higher F1 score than Chan and Glove-WK, but not significantly higher than Glove-Twitter. The results suggest that the UD embeddings, along with Word2Vec and Glove-Twitter, place offensive words semantically close to other words from the same Hurtlex categories, indicating that these embeddings better reflect the categorisation of terms outlined in Hurtlex.

Additionally, we hypothesised that (a) Word2Vec will perform the best at classifying offensive words that are related to minorities, which are in the PS, IS, RE, QAS, and CDS classes, (b) Glove-WK will perform the best for words related to homosexuality, which are in the OM, and CDS classes, and (c) Glove-Twitter, UD, and Chan will perform best for words related to women, which are in ASF, OM, PR, and CDS classes. The results showed that our hypothesis



Figure 3: F1 scores for each class of the kNN model using each word embedding on the Hurtlext test set

holds for UD regarding OM, ASF, and PR and for Word2Vec regarding RE and QAS. However, for the rest of the word embeddings, our hypotheses do not hold, as Glove-Twitter and Glove-WK perform the best at classifying the words in the IS category, where Word2Vec was expected to perform the best, while Chan did not outperform any other word embeddings. Consequently, the acquired results do not provide conclusive answers to how the SOS bias in word embeddings influences the downstream task of offensive words categorisation.

A.2 Racial bias

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1013

1014

1015

1016

1017

1018 1019

1020

1021

1023

1025

1026

1028

1029

1030

1031 1032

1033

1034

1035

To measure the racial bias using the state-of-the-art metrics, we used two target groups: Target group 1, which contains white people's names, and Target group 2, which contains African, Hispanic, and Asian names, and two attribute lists: Attribute list 1, which contains white people occupation names; and Attribute list 2, which contains African, Hispanic, and Asian people's occupations (Badilla et al., 2020; Garg et al., 2018). Then, we measured the average racial bias scores across the different attribute lists for each word embedding using the different metrics (WEAT, RND, RNSB, and ECT). For the SOS bias, we used the mean SOS scores of the words that belong to the "Other ethnicities" category, as computed in Section 3.2 (Figure 1). Finally, we ranked the bias scores as described in Section 5 and computed the Spearman's rank correlation coefficient between the racial bias scores of the different word embeddings and the F1 scores achieved by the two deep learning models on the Twitter-racism and HateEval datasets using the different word embeddings.

The results in Table 9 show that for Twitterracism, SOS has the highest positive correlation with the F1 scores of the MLP model compared to the rest of the bias metrics, whereas WEAT has the highest correlation with the F1 scores of the BiLSTM model. For HateEval, SOS has the highest positive correlation with the F1-scores of the BiLSTM model compared to the rest of the bias metrics, whereas RNSB has the highest correlation with the F1 scores of the MLP model, with SOS only having a higher correlation than WEAT. 1036

1037

1038

1039

1040

1041

1042

Dataset	Model	Spearman's correlation				
	Widdei -	WEAT	RNSB	RND	ECT	SOS
Twitter-racism	MLP	0.200	-0.900	-0.700	-0.200	0.300
	BiLSTM	0.600	-0.700	-0.100	0.100	-0.100
HateEval	MLP	-0.200	0.900	0.300	0.500	0.300
	BiLSTM	-0.205	0.153	-0.718	0.359	0.872

Table 9: Spearman's rank correlation coefficient of the racial bias scores of the different word embeddings and the F1 scores of the deep learning models for each bias metric and dataset.