

An overview of word-embedding methodologies to extinguish bias in deep-learning approaches to toxic comment classification

Hassan R. S. Andrabi

Abstract

Interactions occurring online are typically insensitive to the threat of social accountability for toxic behaviour, leading to frequent employment of abusive and anti-social tactics that go far beyond what might be considered acceptable in face-to-face settings. While applications of machine learning techniques to automatically detect and classify instances of toxic behaviour are well-studied, it has proven difficult to occlude biases in training datasets from flowing onwards to the classifier, and thereafter contributing to discriminatory classifications against sensitive classes such as race, religion, and gender. In this note, I seek to illuminate the capacity for word-embedding techniques to suppress undue learning of sensitive biases in training datasets. I consider this objective in the context of two popular deep-learning frameworks for toxic comment classification: long short-term memory (LSTM) networks, and convolutional neural networks (CNNs). I demonstrate that [...].

1 Introduction

Ours is the age of cyberspace: now, more than ever before, individuals possess an unrestrained freedom to express their opinions for all to behold. This new setting does not reprimand us for opinions expressed anarchically and without requisition; instead, it promotes extreme behaviours that are insular to empathetic considerations. It is against this contextual background that the majority of modern social interactions transpire — social interactions that exploit, by and large, the traceable anonymity of online systems and the resulting protection against social accountability. Indeed, the rise of the internet has transformed the primary setting of our social interactions: but can we be trusted to handle our newfound freedom in a responsible way? In the absence of naturally civilised interactions, moderation becomes crucial to maintaining positive and healthy discussions — however, the quantity and frequency of online communication is simply too vast to be effectively moderated in a manual fashion. This raises an important and imminent question: how do we effectively and automatically monitor the inexorable flow of online interactions to limit toxic behaviour?

Given these analytical problems and the impending social importance of developing measures to promote healthy online interactions, the analysis in this note will empirically

examine the mitigating capacity of machine-learning techniques to identify and filter out textual instances of toxic behaviour. To this end, I focus my analysis on the assessment of the joint classification capacity of: (1) a range popular of machine-learning techniques to accurately classify toxic behaviour; and, (2) a suite of text-embedding representations that propose alternative methods for numerically encoding textual data. In particular, I consider the capacity of these model-embedding representations to minimise incidence of biased classifications: that is, the erroneous tendency for discriminatory classification against sensitive classes such as race, religion, and gender.

The remainder of the paper is structured as follows. Section 2 provides a contextual overview of relevant work in the field of toxic comment classification. Section 3 outlines the dataset and pre-processing methods employed in this analysis. Section 4 describes the general experimental framework, and the architecture of studied classification models. Section 5 presents and discusses the results. Finally, Section 6 concludes the paper

2 Related Work

Prior research has already examined the capacity of various machine-learning and data-engineering techniques to classify toxic comments. Of general interest to this paper is the literature on sentiment analysis, which combines natural-language-processing (NLP) techniques and opinion mining to emulate human-level comprehension of positive or negative sentiment expressed in textual statements [8, 9]. A relatively new branch of NLP literature considers applications of sentiment-analysis to the task of toxic behaviour. Research in this domain can be stratified with respect to the specific dimension of toxic behaviour of interest: besides general classification of toxic online comments [14, 25, 29], related literature also considers classification of specific dimensions of toxic behaviour, including hate speech [4, 19, 26, 31]; harassment [1, 6, 18]; abusive-language [7, 20, 30]; and cyber-bullying [2, 11, 16].

With respect to the selection of machine-learning frameworks employed for toxic comment classification, prior research is generally consistent in its advocacy of certain machine learning techniques as better-suited for the task of toxic comment classifications. In particular, approaches towards classification of toxicity appear to prefer employment of convolutional neural networks (CNNs) [3], as these methods are known to excel at tasks involving elements of pattern-recognition. While CNNs are perhaps most well-known for applications in image-recognition [22], effective translation to sentiment-analysis tasks is unsurprising given the role of syntactical pattern recognition in language comprehension. For example, syntactical patterns such as word order, indicative phrases, and idioms all modify meaning in a way that is algorithmic and theoretically 'learnable'. Indeed, understanding the implications of such patterns is essential to accurate comprehension of language.

At the data pre-processing level, a number of studies consider potential for improved toxic-comment classification through pre-processing using word-embedding techniques, such as TF-IDF [5, 10, 27], GloVe [5, 12, 24], Word2Vec [13], and FastText [24] (for a recent review of these techniques, see Birunda and Devi, 2021 [28]). These techniques systematically estimate vector representations for words in a specified vocabulary, such that words arising from common contexts exhibit similar vector representations. In general, most research in

this field is hampered by challenges associated with skewed class distributions, leading to uneven training exposure to different classes of toxicity. To some extent, these challenges can be mitigated by intentional downsampling of training datasets to impose equal class distributions — although this comes at the cost of less overall data for model training.

3 Dataset

The analysis in this note uses the Jigsaw/Conversation AI Unintended Bias in Toxicity Classification competition dataset (available online: <https://www.kaggle.com/competitions/jigsaw-unintended-bias-in-toxicity-classification/overview/description>). The dataset contains 155,000 annotated comments collected from an archive of the Civil Comments platform: a commenting plugin for online news sites. These comments were annotated by human raters using binary toxicity labels, as well as a series of binary identity labels representing social identities mentioned in the comments. To obtain toxicity labels, each comment was presented to at least 10 human raters, who were prompted to rate comment toxicity according to predefined criteria presented in Table 1. Notably, all comments included in this dataset were subject to a peer-review screening process imposed by Civil Comments. This manual peer-review system was designed to filter out obvious instances toxicity, and substantially limits diversity of vocabulary across the dataset. In particular, the dataset contains very few instances of profane language, and is unlikely to generalise effectively to contexts with less restrictive tenets of commenting etiquette. Figures 1 and 2 present visualisations of the most frequent words appearing in toxic and non-toxic comments respectively.

3.1. Pre-processing

Textual content of comments were cleaned and pre-processed using word-embedding techniques. In particular, all comments were normalised to lower-case, stripped of punctuation marks and non-alphabetic characters, and then tokenised into vector representations through word-embedding techniques. I apply pre-processing using three popular word-embedding representations: term frequency - inverse document frequency (TF-IDF) [15, 17]; Global Vectors for Word Representation (GloVe) [21]; and Sentence-BERT [23]. Each of these techniques attempts to generate vector representations of textual content, such that sentences arising from similar contexts exhibit similar vector representations. Prior to GloVe based embedding, comments were padded or truncated to a token-length of 100, in order to ensure consistent dimensionality of inputs required for model estimation. Thereafter, the total dataset is partitioned into train and validation sets, with 140,000 and 15,000 instances allocated to each set respectively.

4 Experimental method

The following section outlines model architectures and estimation methodology employed to extract dependency between words and phrases in textual comments. To explicate further, consider the following example comment: “Muslims hate gays and want them dead”.

When viewed in isolation, individual words such as: “Muslims”, and “gays” are not particularly indicative of toxic motivations — these words may appear in a variety of perfectly healthy discussions. Toxic motives online become discernible when words appear in particular combinations, such as “hate gays”, or “want them dead”. In such examples, it is clear that effective classification of intent requires an understanding of the encoding of syntactic patterns based on relative positions of critical words. Indeed, these are the patterns that word-embedding methodologies aim to capture.

Leveraging three independent word-embedding techniques introduced in Section 3, I estimate the classification performance of three popular machine-learning frameworks: k-nearest neighbours; (2) logistic regression; and (3) Convolutional Neural network. To evaluate the models, I employ a series of standard metrics used in classification tasks: accuracy, precision, recall, and F1 score. I assess model performance using the total dataset, and then separately across subsets of comments mentioning particular identity subclasses. In total, my analysis estimates fourteen combinations of model structure and word-embedding, summarised in Table 2.

4.1. Logistic regression model

The analysis implements a binary logistic regression model to predict toxicity. To account for proportionally low prevalence of ‘toxic’ labelled comments in the input dataset, penalties for toxicity class weights are set to be inversely proportional to the prevalence of classes in the input dataset. Subsequently, the estimation process imposes larger penalties for inaccurate classifications of ‘toxic’ labels, as compared to ‘non-toxic’ labels. Estimation occurs with L2 regularisation of weights.

4.2. Shallow neural network (NN) model

Given substantial non-linear processing has already occurred in pre-trained sentence-BERT and GloVe word-embeddings, an extensive number of hidden layers may not be required to sufficiently capture unaccounted non-linearity in the input-output relationship. Accordingly, I implement a shallow neural network (NN) consisting of a pre-trained word-embedding layer, a single fully-connected hidden layer with 128 nodes, and an output layer with a sigmoid activation function. Hidden nodes are configured with Rectified Linear-unit (ReLU) activation functions.

4.3. Deep neural network (DNN) model

Deep neural networks are well-studied to exhibit exceptional performance in sentiment-analysis tasks. I implement a DNN consisting of a pre-trained word-embedding layer, 2 fully-connected hidden layers single hidden layer with 128 nodes each, and an output layer with a sigmoid activation function. As with the NN, hidden nodes are configured with Rectified Linear-unit (ReLU) activation functions.

4.4. Convolutional neural network (CNN) model

Among the considered approaches, applications of convolutional neural networks to the sentiment-classification tasks are perhaps the most well-studied. I estimate a convolutional neural network (CNN) consisting of a pre-trained word-embedding layer, a single fully-connected hidden layer with 128 nodes, and an output layer with a sigmoid activation function. Hidden nodes are configured with Rectified Linear-unit (ReLU) activation functions.

5 Results

Model estimation and evaluation is conducted in congruence with hold-out methodology, with a train-test split of $N=140,000$, and $N=15,000$ respectively. Each model, in turn, was separately trained on the training dataset using TF-IDF, GloVe, and sentence-BERT word-embedding methodologies, and classification performance was evaluated on a consistently embedded validation dataset. Performance is assessed by way of precision, recall, and F-score metrics between toxicity predictions passed by the model, and toxicity annotations labelled by human raters, using the criteria defined in Table 1.

6 Conclusion

References

- [1] S Abarna, JI Sheeba, S Jayasrilakshmi, and S Pradeep Devaneyan. Identification of cyber harassment and intention of target users on social media platforms. *Engineering applications of artificial intelligence*, 115:105283, 2022.
- [2] Arnisha Akhter, Uzzal K Acharjee, and Md Masbaul A Polash. Cyber bullying detection and classification using multinomial naïve bayes and fuzzy logic. *Int. J. Math. Sci. Comput*, 5(4):1–12, 2019.
- [3] Darko Androcec. Machine learning methods for toxic comment classification: a systematic review. *Acta Universitatis Sapientiae, Informatica*, 12(2):205–216, 2020.
- [4] Femi Emmanuel Ayo, Olusegun Folorunso, Friday Thomas Ibharalu, and Idowu Ademola Osinuga. Machine learning techniques for hate speech classification of twitter data: State-of-the-art, future challenges and research directions. *Computer Science Review*, 38:100311, 2020.
- [5] Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. Deep learning for hate speech detection in tweets. In *Proceedings of the 26th international conference on World Wide Web companion*, pages 759–760, 2017.
- [6] Priyam Basu, Tiasa Singha Roy, Soham Tiwari, and Saksham Mehta. Cyberpolice: Classification of cyber sexual harassment. In *EPIA Conference on Artificial Intelligence*, pages 701–714. Springer, 2021.
- [7] Peter Bourgonje, Julian Moreno-Schneider, Ankit Srivastava, and Georg Rehm. Automatic classification of abusive language and personal attacks in various forms of online communication. In *International Conference of the German Society for Computational Linguistics and Language Technology*, pages 180–191. Springer, Cham, 2017.
- [8] Erik Cambria and Bebo White. Jumping nlp curves: A review of natural language processing research. *IEEE Computational intelligence magazine*, 9(2):48–57, 2014.
- [9] KR1442 Chowdhary. Natural language processing. *Fundamentals of artificial intelligence*, pages 603–649, 2020.
- [10] Claudio Moisés Valiense de Andrade and Marcos André Gonçalves. Profiling hate speech spreaders on twitter: Exploiting textual analysis of tweets and combinations of multiple textual representations. In *CEUR Workshop Proc*, volume 2936, pages 2186–2192, 2021.
- [11] Michele Di Capua, Emanuel Di Nardo, and Alfredo Petrosino. Unsupervised cyber bullying detection in social networks. In *2016 23rd International conference on pattern recognition (ICPR)*, pages 432–437. IEEE, 2016.
- [12] Antigoni Maria Founta, Despoina Chatzakou, Nicolas Kourtellis, Jeremy Blackburn, Athena Vakali, and Ilias Leontiadis. A unified deep learning architecture for abuse detection. In *Proceedings of the 10th ACM conference on web science*, pages 105–114, 2019.

- [13] Björn Gambäck and Utpal Kumar Sikdar. Using convolutional neural networks to classify hate-speech. In *Proceedings of the first workshop on abusive language online*, pages 85–90, 2017.
- [14] Spiros V Georgakopoulos, Sotiris K Tasoulis, Aristidis G Vrahatis, and Vassilis P Plagianakos. Convolutional neural networks for toxic comment classification. In *Proceedings of the 10th hellenic conference on artificial intelligence*, pages 1–6, 2018.
- [15] Karen Sparck Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 1972.
- [16] Tarek Kanan, Amal Aldaaja, and Bilal Hawashin. Cyber-bullying and cyber-harassment detection using supervised machine learning techniques in arabic social media contents. *Journal of Internet Technology*, 21(5):1409–1421, 2020.
- [17] Hans Peter Luhn. A statistical approach to mechanized encoding and searching of literary information. *IBM Journal of research and development*, 1(4):309–317, 1957.
- [18] Tolba Marwa, Ouadfel Salima, and Meshoul Souham. Deep learning for online harassment detection in tweets. In *2018 3rd International Conference on Pattern Analysis and Intelligent Systems (PAIS)*, pages 1–5. IEEE, 2018.
- [19] Nanlir Sallau Mullah and Wan Mohd Nazmee Wan Zainon. Advances in machine learning algorithms for hate speech detection in social media: a review. *IEEE Access*, 2021.
- [20] Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. Abusive language detection in online user content. In *Proceedings of the 25th international conference on world wide web*, pages 145–153, 2016.
- [21] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [22] Waseem Rawat and Zenghui Wang. Deep convolutional neural networks for image classification: A comprehensive review. *Neural computation*, 29(9):2352–2449, 2017.
- [23] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- [24] Alison Ribeiro and Nádia Silva. Inf-hateval at semeval-2019 task 5: Convolutional neural networks for hate speech detection against women and immigrants on twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 420–425, 2019.
- [25] Julian Risch and Ralf Krestel. Toxic comment detection in online discussions. In *Deep learning-based approaches for sentiment analysis*, pages 85–109. Springer, 2020.
- [26] Georgios Rizos, Konstantin Hemker, and Björn Schuller. Augment to prevent: short-text data augmentation in deep learning for hate-speech classification. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pages 991–1000, 2019.

- [27] Punyajoy Saha, Binny Mathew, Pawan Goyal, and Animesh Mukherjee. Hateminers: Detecting hate speech against women. *arXiv preprint arXiv:1812.06700*, 2018.
- [28] S Selva Birunda and R Kanniga Devi. A review on word embedding techniques for text classification. *Innovative Data Communication Technologies and Application*, pages 267–281, 2021.
- [29] Betty Van Aken, Julian Risch, Ralf Krestel, and Alexander Löser. Challenges for toxic comment classification: An in-depth error analysis. *arXiv preprint arXiv:1809.07572*, 2018.
- [30] Bertie Vidgen and Leon Derczynski. Directions in abusive language training data, a systematic review: Garbage in, garbage out. *Plos one*, 15(12):e0243300, 2020.
- [31] Fan Yang, Xiaochang Peng, Gargi Ghosh, Reshef Shilon, Hao Ma, Eider Moore, and Goran Predovic. Exploring deep multimodal fusion of text and photo for hate speech classification. In *Proceedings of the third workshop on abusive language online*, pages 11–18, 2019.

List of Figures

- 1 Words appearing most frequently in toxic comments. Frequencies represent the number of times an indicated word appears in a comment with a 'toxic' annotation. The top twenty-five most frequent words are labelled. 10
- 2 Words appearing most frequently in non-toxic comments. Frequencies represent the number of times an indicated word appears in a comment with a 'non-toxic' annotation. The top twenty-five most frequent words are labelled. 10

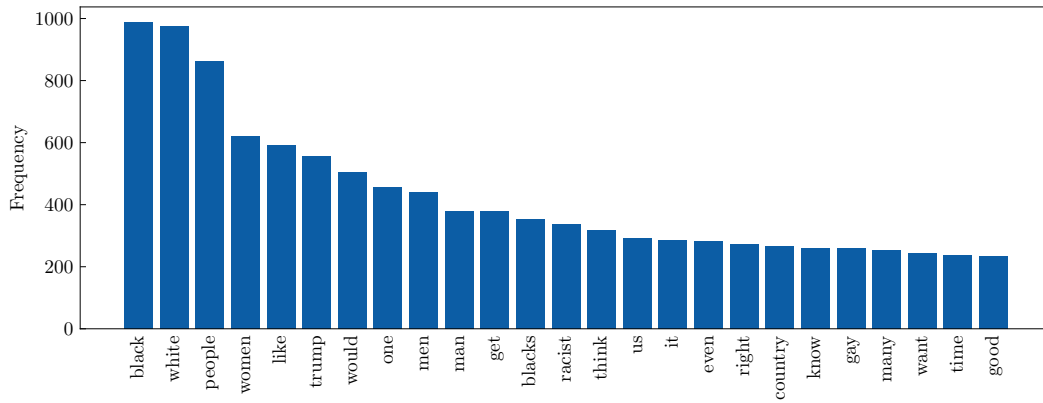


Figure 1: Words appearing most frequently in toxic comments. Frequencies represent the number of times an indicated word appears in a comment with a 'toxic' annotation. The top twenty-five most frequent words are labelled.

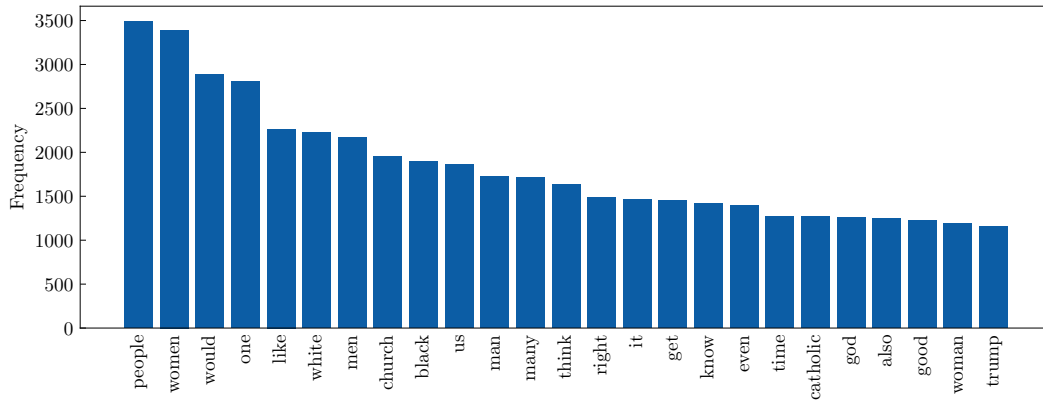


Figure 2: Words appearing most frequently in non-toxic comments. Frequencies represent the number of times an indicated word appears in a comment with a 'non-toxic' annotation. The top twenty-five most frequent words are labelled.

List of Tables

1	Jigsaw/Coversation AI toxicity labelling criteria	12
2	Model structure and word-embedding combinations	12
3	Toxic comment classification performance	13

Table 1: Jigsaw/Coversation AI toxicity labelling criteria

Label	Criteria
Very Toxic	A very hateful, aggressive, or disrespectful comment that is very likely to make you leave a discussion or give up on sharing your perspective
Toxic	A rude, disrespectful, or unreasonable comment that is somewhat likely to make you leave a discussion or give up on sharing your perspective
Hard to say	No criteria given
Not toxic	No criteria given

Table 2: Model structure and word-embedding combinations

Model	Word-embedding methodology
Logistic regression	TF-IDF embeddings
Logistic regression	sentence-BERT embeddings
Logistic regression	GloVe embeddings (minimum across 100 dimensions)
Logistic regression	GloVe embeddings (maximum across 100 dimensions)
Logistic regression	GloVe embeddings (average across 100 dimensions)
Shallow neural network (NN)	TF-IDF embeddings
Shallow neural network (NN)	sentence-BERT embeddings
Shallow neural network (NN)	GloVe embeddings
Deep neural network (DNN)	TF-IDF embeddings
Deep neural network (DNN)	sentence-BERT embeddings
Deep neural network (DNN)	GloVe embeddings
Convolutional neural network (CNN)	TF-IDF embeddings
Convolutional neural network (CNN)	sentence-BERT embeddings
Convolutional neural network (CNN)	GloVe embeddings

Table 3: Toxic comment classification performance

Model	Precision	Accuracy	Recall	F-score
logit (TF-IDF)	0.6376	0.8276	0.2035	0.3086
logit (sentence-BERT)	0.6506	0.8335	0.2575	0.3690
logit (GloVe - min.)	0.5718	0.8143	0.0702	0.1250
logit (GloVe - max.)	0.5915	0.8165	0.0935	0.1614
logit (GloVe - avg.)	0.5833	0.8170	0.1111	0.1867
NN (TF-IDF)	0.6031	0.8215	0.1630	0.2566
NN (sentence-BERT)	0.6040	0.8254	0.2212	0.3238
NN (GloVe)	0.6438	0.8397	0.3404	0.4453
DNN (TF-IDF)	0.0000	0.8110	0.0000	0.0000
DNN (sentence-BERT)	0.6341	0.8295	0.2310	0.3387
DNN (GloVe)	0.7043	0.8380	0.2462	0.3649
CNN (TF-IDF)	0.6902	0.8220	0.3034	0.4194
CNN (sentence-BERT)	0.6983	0.8326	0.3199	0.4381
CNN (GloVe)	0.7011	0.8460	0.3228	0.4420