

Problem statement

Fair treatment by a BERT-based Twitter disinformation classifier

We use a quantitative bias scan tool to assess fair treatment of a self-trained disinformation detection algorithm on Twitter data. This document presents statistically significant disparities found by the tool. The results are submitted to a commission of human experts. This audit commission formulates normative advice if, and how, (higher-dimensional) proxy discrimination and/or ethically undesirable forms of differentiation could be investigated further.

1. Introduction

Unfair treatment by algorithms is multi-faceted. A first concern is one-dimensional proxy discrimination. Proxy discrimination concerns unlawful differentiation based on an apparently neutral feature (such as *literacy rate*) that is critically linked to a protected ground as specified in legal directives¹ (such as *ethnicity*). A second concern is ethically undesirable forms of differentiation. Algorithms can differentiate upon a seemingly innocuous feature, such as browser type or house number suffix. This type of differentiation evades non-discrimination law, as many features are not critically linked to a protected ground, but can still be perceived as unfair, for instance if it reinforces socio-economic inequality. A third concern is higher-dimensional forms of unfair treatment. Algorithms differentiate upon clusters that are defined by a mixture of features. Higher-dimensional forms of algorithmic differentiation are difficult to detect for humans. Let alone to assess whether the cluster is involved in proxy discrimination and/or ethically undesirable forms of differentiation. In theory, statistical methods are capable to detect both higher- and one-dimensional forms of undesirable differentiation. In this case study, we use a statistical bias scan tool to examine in practice whether the above concerns can be overcome.

¹ In the European Union (EU), the European Convention of Human Rights (ECHR) serves as the legal fundament against discrimination. Additional EU directives (2000/43/EC, 2000/78/EC, 2004/113/EC, and 2006/54/EC) provide context-specific protection, e.g., persons with disabilities, employment rights, and consumer protection.

2. Unsupervised bias scan

The bias scan tool² identifies clusters for which a binary classification algorithm is systematically misclassifying, i.e., predicting a different class than the ground truth label in the data. A cluster is a group of datapoints sharing similar features. The tool makes use of unsupervised clustering³ and therefore does not require *a priori* information about existing disparities and protected attributes of users (which are often not available in practice).

For this case study, we review a BERT-based disinformation classification algorithm⁴ which is trained on the Twitter1516 dataset⁵, enriched with self-collected Twitter API data⁶. The dataset consists of 1,057 verified true and false tweets, 3 user features (verified profile, #followers, user engagement) and 5 content features (length, #URLs, #mentions, #hashtags, sentiment score). We run two bias scans. In Scan 1, the bias metric is defined by the False Positive Rate (FPR). FPR relates to true content predicted to be false, proportional to all true content. In Scan 2, the bias metric is defined by the False Negative Rate (FNR). FNR relates to false content predicted to be true, proportional to all false content. In sum:

Scan 1. Bias = $FPR_{\text{cluster}} - FPR_{\text{rest of dataset}}$

Scan 2. Bias = $FNR_{\text{cluster}} - FNR_{\text{rest of dataset}}$

The full bias scan pipeline is displayed in Figure 1.



Figure 1 – Bias scan pipeline for the disinformation classifier.

² Misztal-Radecka, Indurkya, Bias-Aware Hierarchical Clustering for detecting the discriminated groups of users in recommendation systems, *Information Processing and Management* (2021).

³ Documentation about the k-means Hierarchical Bias-Aware Clustering (HBAC) algorithm:
https://github.com/NGO-Algorithm-Audit/Bias_scan/blob/master/Technical_documentation_bias_scan.pdf

⁴ More information about the self-trained BERT-based classification algorithm:
https://github.com/NGO-Algorithm-Audit/Bias_scan/blob/master/HBAC_scan/HBAC_BERT_disinformation_classifier.ipynb

⁵ Liu, Xiaomo and Nourbakhsh, Armineh and Li, Quanzhi and Fang, Rui and Shah, Sameena, in *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management* (2015)

⁶ More information on the data collection process:
https://github.com/NGO-Algorithm-Audit/Bias_scan/blob/master/data/Twitter_dataset/Twitter_API_data_collection.ipynb

3. Results: Identified quantitative disparities

For Scan 1, the cluster for which the disinformation classifier is underperforming the most (bias=0.08, n=249) is characterized by the features displayed in Figure 2.

Difference in means is the difference in standardized feature means between the disparately treated cluster and the rest of the dataset. Hypothesis testing⁷ indicates that on average, user that:

- are verified, have higher #followers, user engagement and #URLs;
- use less #hashtags and have lower tweet length

have more true content classified as false (false positives).

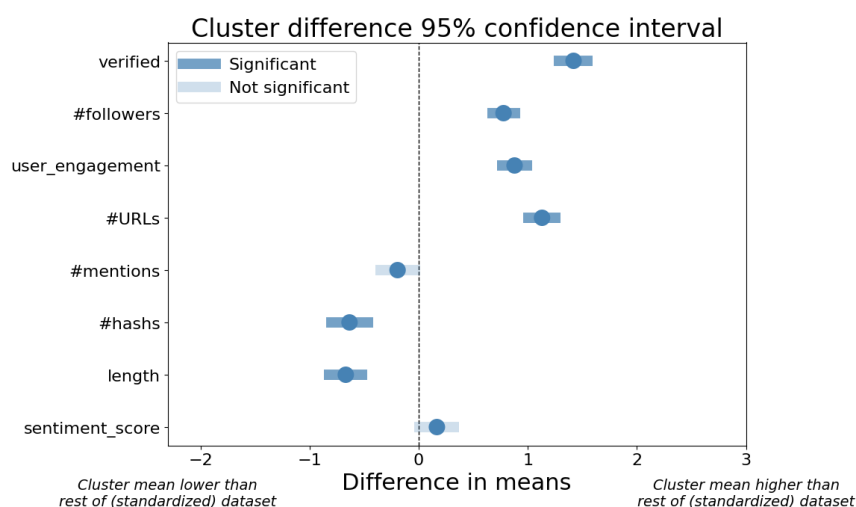


Figure 2 – Identified quantitative feature disparities in cluster with highest bias.
For this bias scan, bias is defined by the False Positive Rate.

For Scan 2, the cluster for which the disinformation classifier is underperforming the most (bias=0.13, n=46) is characterized by the features displayed in Figure 3.

Hypothesis testing indicates that on average, user that:

- use more #hashtags and have higher sentiment score;
- are non-verified, have less #followers, user engagement and tweet length

have more false content classified as true (false negatives).

⁷ Here, the hypothesis tested is that there is no difference in feature means of the cluster and the pooled feature means of other clusters. These differences are statistically significant even after performing a Bonferroni correction to adjust for false discoveries due to multiple hypothesis testing.

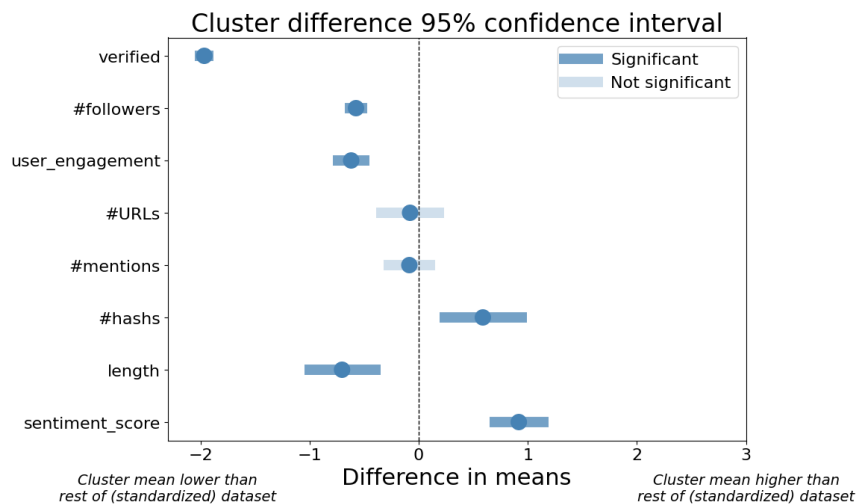


Figure 3 – Identified quantitative feature disparities in cluster with highest bias. For this bias scan, bias is defined by the False Negative Rate.

These results might indicate (higher-dimensional) unfair treatment by the disinformation classifier. More information on the identified clusters and robustness tests of the results can be found in the Appendix.

4. Qualitative assessment of identified disparities

The identified disparities in Section 3 do not establish prohibited *prima facie* discrimination. Rather, the identified disparities serve as a starting point to assess potential unfair treatment according to the context-sensitive qualitative doctrine. To assess unfair treatment, we question:

- i) Is there an indication that one of the statistically significant features, or a combination of the features, stated in Figure 2-3 are critically linked to one or multiple protected grounds?
- ii) In the context of disinformation detection, is it as harmful to classify true content as false (false positive) as false content as true (false negative)?
- iii) For a specific cluster of people, is it justifiable to have true content classified as false 8 percentage points more often? For a specific cluster of people, is it justifiable to have false content classified as true 13 percentage points more often?
- iv) Is it justifiable that the disinformation classification algorithm is too harsh towards users with verified profile, more #followers and higher user engagement and too lenient towards users with non-verified profile, less #followers and lower user engagement?

Auditing disinformation detection algorithms

As of December 2022, Article 28 of the European Digital Services Act (DSA) subjects very large online platforms (VLOPs) to annual independent auditing of their services and risk mitigation measures. Open-source AI auditing tools, such as this bias scan tool, help to detect and mitigate (higher-dimensional) forms of unfair treatment in disinformation detection and other AI (ranking and recommender) systems.

With this case study, Algorithm Audit aims to provide qualitative guidelines how statistical methods can be used to monitor unfair treatment by AI systems. Without clear guidance from data protection boards, supervisors, researchers and algorithmic regulatory bodies, misinterpretation of quantitative metrics stands in the way of independent quantitative and qualitative oversight of the risk of biased AI systems. Building on the quantitative results of the bias scan, Algorithm Audit provides qualitative justifications to make a normative judgment about whether AI systems are causing unfair treatment or not.