

# Classifying Fake News on the Fake and Real news Dataset

Pedro Antônio Silva  
Ciência da Computação

INSPER  
Instituto de Ensino e Pesquisa

SP, São Paulo  
pedroas5@al.insper.edu.br

## I. DATASET

The dataset used was initially discussed in the context of bridging the gap between AI and human linguistic judgment, as highlighted in the paper "AI vs linguistic-based human judgement: Bridging the gap in pursuit of truth for fake news detection" [1]. The business application, is therefore, to leverage AI models to assist human fact-checkers in quickly filtering out potentially harmful or misleading content.

## II. CLASSIFICATION PIPELINE

The classification pipeline began with text pre-processing using Lemmatization [2], to reduce words to their root form, then a TF-IDF (Term Frequency-Inverse Document Frequency) [3] vectorizer was applied to extract features from the text to highlight the importance of word frequency in distinguishing true and false news articles.

A Logistic Regression classifier [4] was then applied to the TF-IDF vectors. this model was used due to its interpretability and simplicity, where the learned coefficients indicate the weight given to each word in the classification process. The model assumes that certain terms will be highly indicative of fake news, such as emotionally charged or sensational words, internet lingo and colloquial terms.

## III. EVALUATION

A balanced accuracy score was used to evaluate the classifier, ensuring the results weren't biased by the imbalance between true and false articles.

The classifier achieved a 94% balanced accuracy score, which aligns with findings from the paper Fake News Detection Platform—Conceptual Architecture and Prototype [5], cited in the previously referenced work [1]. Significant words from each class were chosen as key examples of what the model considers during classification, based on the learned coefficients:

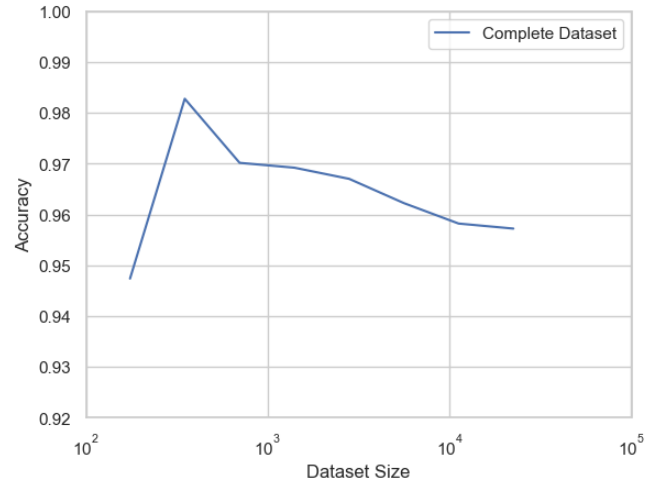
Score	Word	Lemma	Specific Term
-3.69	watch	O	X
-3.13	https	X	O
7.32	democratic	X	X
7.6	edt	X	O
-4.1	read	O	X

<sup>a</sup>. most impactful words for classification

Negative scores, such as for the verbs "watch" and "read," indicate a higher likelihood of classifying an article as fake when these words are present, possibly due to their common use in online media. On the other hand, words like "democratic" had higher positive scores, as their use is more common in legitimate articles in the dataset. Additionally, terms related to media types and goals, such as "edt" (a time zone) and "https" (a link component), had scores reflecting where they commonly occur, with the first contributing positively and the second negatively

## IV. DATA SIZE

As the dataset used has a considerable size and the model achieved high accuracy, a down sampling strategy was applied to seek understanding of the impacts and relevance the amount of data collected had in the model.

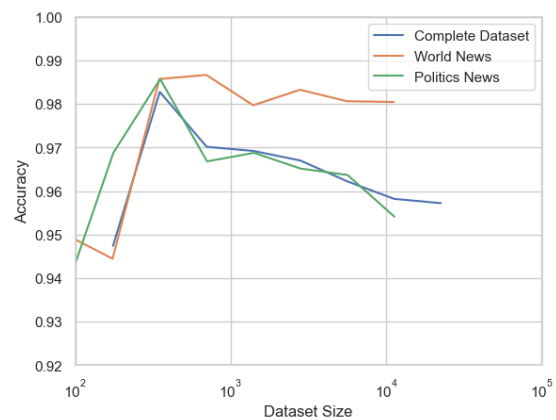


Observing the results, it can be assessed that this model displays a large increase in accuracy as the sample size approaches  $10^3$  and stabilizes while approaching  $10^4$ , with a noticeable diminished gain from sample size increases.

## V. TOPIC ANALYSIS

The analysis began by incorporating NMF (Non-negative Matrix Factorization) [6] between the vectorizer and classifier in the pipeline. After fine-tuning the model to determine the optimal number of topics, 10 topics were selected, this two-step configuration yielded a balanced accuracy score of approximately 88%.

To understand if this loss in accuracy could be corrected through a more refined topic delimitation, the dataset was split into two subjects already defined in the news and the data down sampling experiment re-conducted:



As a conclusion, the collected data indicates that for certain topics a specialized model may yield a better accuracy, but it is not always an improvement, and automated topic classification may not be the most reliable.

## REFERENCES

- [1] AI vs linguistic-based human judgement: Bridging the gap in pursuit of truth for fake news detection..
- [2] **Bird, S., Klein, E., & Loper, E. (2009).** *Natural Language Processing with Python*. O'Reilly Media, Inc.
- [3] **Ramos, J. (2003).** *Using TF-IDF to determine word relevance in document queries*. In *Proceedings of the First International Conference on Machine Learning* (pp. 133-142).
- [4] **D. J. C. MacKay**, Information Theory, Inference, and Learning Algorithms, Cambridge University Press, 2003.
- [5] **R. Kozik, M. Pawlicki, S. Kula, M. Choraś**, Fake news detection platform—conceptual architecture and prototype
- [6] **Lee, D. D., & Seung, H. S. (1999).** "Learning the parts of objects by non-negative matrix factorization." *Nature*, 401(6755), 788–791. doi:10.1038/4456