

Hoaxpedia: A Unified Wikipedia Hoax Articles Dataset

Hsuvas Borkakoty¹, Luis Espinosa-Anke^{1,2},

¹Cardiff NLP, School of Computer Science and Informatics, Cardiff University, UK

²AMPLYFI, UK

{borkakotyh,espinosaankel}@cardiff.ac.uk

Abstract

Hoaxes are a recognised form of disinformation created deliberately, with potential serious implications in the credibility of reference knowledge resources such as Wikipedia. What makes detecting Wikipedia hoaxes hard is that they often are written according to the official style guidelines. In this work, we first provide a systematic analysis of the similarities and discrepancies between legitimate and hoax Wikipedia articles, and introduce HOAXPEDIA, a collection of 311 Hoax articles (from existing literature as well as official Wikipedia lists) alongside semantically similar real articles. We report results of binary classification experiments in the task of predicting whether a Wikipedia article is real or hoax, and analyze several settings as well as a range of language models. Our results suggest that detecting deceitful content in Wikipedia based on content alone, despite not having been explored much in the past, is a promising direction.¹

1 Introduction

Wikipedia is, as Hovy et al. (2013) define it, the “largest and most popular collaborative and multilingual resource of world and linguistic knowledge”, and it is acknowledged that its accuracy is on par with or superior than, e.g., the Encyclopedia Britannica (Giles, 2005). However, as with any other platform, Wikipedia is also the target of online vandalism, and *hoaxes*, a more obscure, less obvious form of vandalism², constitute a threat to the overall integrity of this collaborative encyclopedia (Kumar et al., 2016; Wong et al., 2021; Wang and McKeown, 2010), precisely because of its “publish

first, ask questions later” policy (Asthana and Hal-faker, 2018). Although Wikipedia employs community based New Page Patrol systems to check the credibility of a newly created article, the process is always in backlog³, making the process overwhelming (Schneider et al., 2014).

Hoax articles, created to deliberately spread false information (Kumar et al., 2016), harm the credibility of Wikipedia as a knowledge resource, and generate concerns among its users (Hu et al., 2007). Since manual inspection of quality is a process typically in backlog (Dang and Ignat, 2016), the automatic detection of such articles is a desirable feature. However, most works in the literature have centered their efforts in metadata associated with hoax articles, e.g., user activity, appearance features or revision history (Zeng et al., 2006; Elebiary and Ciampaglia, 2023; Kumar et al., 2016; Wong et al., 2021; Hu et al., 2007; Susuri et al., 2017). For example, Adler et al. (2011) introduced a vandalism detection system using metadata, content and author reputation features, whereas Kumar et al. (2016) provides a comprehensive study of hoax articles and their timeline from discovery to deletion. In this work, the authors define the characteristics of a successful hoax, with a data-driven approach based on studying a dataset of 64 articles (both hoax and real), on top of which they train statistical classifiers. Furthermore, other works have compared network traffic and features of hoax articles to those of other articles published the same day (Elebiary and Ciampaglia, 2023), and conclude that hoax articles attract more attention after creation than *cohort* articles. Finally, Wong et al. (2021) study various Wikipedia vandalism types and introduce the Wiki-Reliability dataset, which comprises articles based on 41 author-compiled templates. This dataset contains 1,300 articles marked as hoax,

¹The Dataset(view-only, access upon request) is available in: https://osf.io/rce8m/?view_only=ed469941644c496fb4a6425297ced1f2. We will publicly release our models and the datasets in Huggingface upon acceptance

²https://en.wikipedia.org/wiki/Wikipedia:Do_not_create_hoaxes.

³https://en.wikipedia.org/wiki/Wikipedia:New_pages_patrol.

which are real articles with false information, a.k.a hoax facts (Kumar et al., 2016).

We part ways from previous works and focus exclusively on the content of hoax articles, and aim to answer the research question “*Can hoax articles be distinguished from real articles using NLP techniques by looking exclusively at the article’s content*”? We first construct a dataset (HOAXPEDIA) containing 311 hoax articles and around 30,000 plausible negative examples, i.e., real Wikipedia articles that are semantically similar to hoax articles, to create a set of negatives that *cover similar topics* to hoax articles (e.g., a newly discovered species). We also explore whether a Wikipedia definition (the first sentence of the article) can provide hints towards its veracity. Our results (reported at different ratios of hoax vs real articles) suggest that, while style and shallow features are certainly not good predictors of hoax vs real Wikipedia articles, LMs are capable of exploiting other more intricate features, and open a promising research direction focused on content-based hoax flagging. Our contributions in this work can be summarised as follows.

- We systematically contrast a set of proven Wikipedia hoax articles vs. legitimate articles.
- We propose HoaxPedia, a novel Wikipedia Hoax article dataset with 311 hoax articles and semantically similar legitimate articles collected from Wikipedia.
- We conduct binary classification experiments on this dataset, using a range of language models to accurately predict whether an article is a hoax or real based solely on its content.

2 HOAXPEDIA Construction

HOAXPEDIA is constructed by unifying five different resources that contain known hoaxes, e.g., from Kumar et al. (2016); Elebiary and Ciampaglia (2023), as well as the official Wikipedia hoaxes list⁴ and the Internet Archive. We used Internet Archive to manually retrieve Wikipedia pages that are now deleted from Wikipedia, but were at one point in the past identified as hoaxes. We manually verify each of the article we collect from Wikipedia and Internet Archive as hoax using their accompanied deletion discussion and reasons for citing them

⁴https://en.wikipedia.org/wiki/Wikipedia:List_of_hoaxes_on_Wikipedia

as a hoax. In terms of negative examples, while we could have randomly sampled Wikipedia pages, this could have introduced a number of biases in the dataset, e.g., hoax articles contain historical events, personalities or artifacts, and thus we are interested in capturing a similar breadth of topics, entities and sectors in the negative examples so that a classifier cannot rely on these spurious features. We select these negative examples (truthful articles) by ensuring they correspond to authentic content. This is achieved by verifying they do not carry the Db-hoax flag, which Wikipedia’s New Page Patrol policy uses to mark potential hoaxes. Within this set, we extract negative examples as follows. Let H be the set of hoax articles, and W the set of candidate *real* Wikipedia pages, with $T_H = \{t_{H^1}, \dots, t_{H^p}\}$ and $T_W = \{t_{W^1}, \dots, t_{W^q}\}$ their corresponding vector representations, and p and q the number of hoax and candidate Wikipedia articles, respectively. Then, for each SBERT(all-MiniLM-L6-v2) (Reimers and Gurevych, 2019) title embedding $t_{H^i} \in T_H$, we retrieve its top k nearest neighbors (NN) from T_W via cosine similarity COS. We experiment with different values for k , specifically $k \in \{2, 10, 100\}$:

$$\text{NN}(t_{H^i}) = \{t_{W^j} : j \in J_k(t_{H^i})\}$$

where $J_k(t_{H^i})$ contains the top k cosine similarities in T_W for a given t_{H^i} , and

$$\text{COS}(t_{H^i}, t_{W^j}) = \frac{t_{H^i} \cdot t_{W^j}}{\|t_{H^i}\| \|t_{W^j}\|}$$

The result of this process is a set of positive (hoax) articles and a set of negative examples we argue will be similar in content and topic, effectively removing any topic bias from the dataset.

3 Hoax vs. Real, a Surface-Level Comparison

To maintain the longevity and to avoid detection, hoax articles follow Wikipedia guidelines and article structure. This raises the following question: “*how (dis)similar are hoaxes with respect to a hypothetical real counterpart?*”. Upon inspection, we found comments in the deletion discussions such as “*I wouldn’t have questioned it had I come across it organically*” (for the hoax article *The Heat is On*⁵), or “*The story may have a “credible feel” to it, but it lacks any sources*”, a comment on article *Chu Chi*

⁵[https://en.wikipedia.org/wiki/Wikipedia:Articles_for_deletion/The_Heat_Is_On_\(TV_series\)](https://en.wikipedia.org/wiki/Wikipedia:Articles_for_deletion/The_Heat_Is_On_(TV_series))

Zui⁶. Comments like these highlight that hoaxes are generally well written (following Wikipedia’s guidelines), and so we proceed to quantify their stylistic differences in a comparative analysis that looks at: (1) article text length comparison; (2) sentence and word length comparison; and (3) a readability analysis.

Article Text length distribution: Following the works of Kumar et al. (2016), we conduct a text length distribution analysis with hoax and real articles, and verify they show a similar pattern (as shown in Figure 1), with similar medians for hoax and real articles, specifically 1,057 and 1,777 words, respectively.

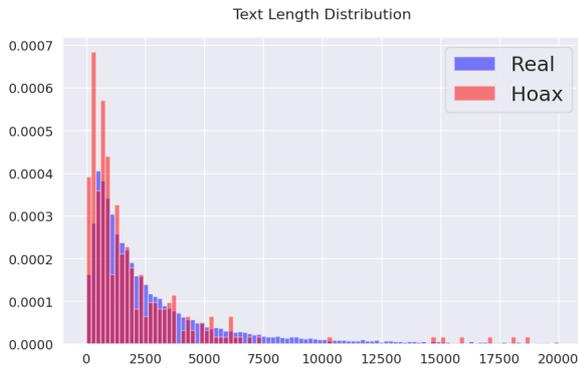


Figure 1: Text length distribution for hoax and real articles.

Average sentence and word length: Calculating average sentence and word length for hoaxes and real articles separately can be a valuable proxy for identifying any obvious stylistic or linguistic (e.g., syntactic complexity) patterns. We visualize these in a series of box plots in Figure 2. They clearly show a similar style, with sentence and word length medians at 21.23 and 22.0, and 4.36 and 4.35 for real and hoax articles respectively.

Readability Analysis: Readability analysis gives a quantifiable measure of the complexities of texts. It can clarify on easement of understanding the text, revealing patterns that can be used to either disguise disinformation like hoaxes or convey clear, factual content. For readability analysis, we use Flesch-Kincaid (FK) Grading system (Flesch, 2007), a metric that indicates the comprehension difficulty when reading a passage in the context of contemporary academic English. This metric gives us

⁶https://en.wikipedia.org/wiki/Wikipedia:Articles_for_deletion/Chu_Chi_Zui

an aggregate of the complexity of documents, their sentences and words, measured by the average number of words per sentence and the average number of syllables per word. After obtaining an average for both hoax and real articles, we visualize these averages again in Figure 2, we find a median of 9.4 for real articles and 9.5 for hoax articles, again highlighting the similarities between these articles.

The above analysis suggests that hoax Wikipedia articles are indeed well disguised, and so in the next section we propose a suite of experiments for hoax detection based on language models, setting an initial set of baselines for this novel dataset.

4 Experiments

We cast the problem of identifying hoax vs. legitimate articles as a binary classification problem, in which we evaluate a suite of LMs, specifically: BERT-family of models (BERT-base and large (Devlin et al., 2019), RoBERTa-base and large (Liu et al., 2019), Albert-base and large (Lan et al., 2019)), as well as T5 (Base and Large) (Raffel et al., 2020) and Longformer (Base) (Beltagy et al., 2020). We use the same training configuration for the BERT-family of models, T5 models and Longformer, and set the generation objective as *Binary classification* for the T5 models. In terms of data size, we consider the three different scenarios outlined in Section 2 (2x, 10x and 100x negative examples). This approach naturally increases the challenge for the classifiers. The details about the data used in different settings are given in Appendix A.

In addition to the three different settings for positive vs negative ratios, we also explore *how much text is actually needed to catch a hoax*, or in other words, *are definition sentences in hoax articles giving something away?* This is explored by running our experiments on the full Wikipedia articles, on one hand, and on the definition (first sentence alone), on the other. This latter setting is interesting from a lexicographic perspective because it helps us understand if the Wikipedia definitions show any pattern that a model could exploit. Moreover, from the practical point of view of building a classifier that could dynamically “patrol” Wikipedia and flag content automatically, a definition-only model would be more interpretable (with reduced ambiguity and focusing on core meaning/properties of the entity) and could have less parameters (handling smaller vocabularies, and compressed knowledge),

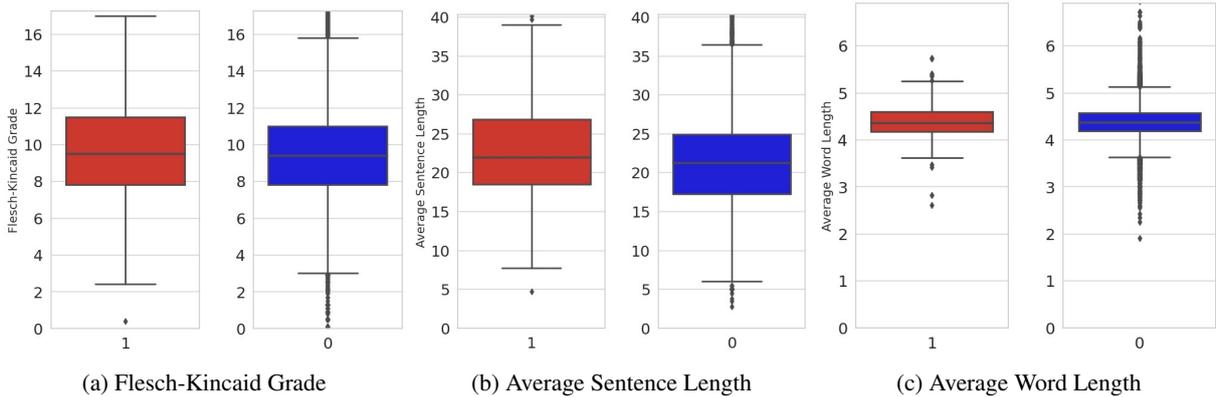


Figure 2: Results of different stylistic analyses on Hoax (red) and Real (blue) articles.

which would have practical retraining/deployment implications in cost and turnaround.

5 Results

Our experiments are aimed to explore the impact of data imbalance and content length. In both cases, we compare several classifiers and analyze whether model size (in number of parameters) is correlated with performance. In terms of evaluation metrics, all results we report are F1 on the positive class (hoax). In definition-only setting, we find that models evaluated on datasets that are relatively balanced (2 Real articles for every hoax) show a stable performance, but they degrade drastically as the imbalance increases. In terms of robustness, RoBERTa is the most consistent, with an F1 of around 0.6 for all three settings, whereas Albert models perform poorly (exhibiting, however, some interesting behaviours, which we will discuss further). For the full text setting, we find that Longformer models performs well, with an F1 of 0.8. Surprisingly, the largest model we evaluated (T5-large) is not the best performing model, although this could point to under-fitting, the dataset potentially being too small for a model this size to train properly. Another interesting behaviour of T5-large is that in the 1 Hoax vs 2 Real setting, it shows exactly the same performance, whether seeing a definition or the full text. On the other side of the spectrum, we find that Albert models are the ones showing the highest improvement when going from definition to full text. This is interesting, as it shows a small model may miss nuances in definitions but can still compete with, or even outperform, larger models.

A perhaps not too surprising observation is that all models improve after being exposed to more

text, as seen in Table 1, increasing their F1 about 20% on average, and sometimes even up to 30. This confirms that definitions alone are not a sufficiently strong signal for detecting hoax articles, although there are notable exceptions. Moreover, in terms of absolute performance, RoBERTa models perform decently, although significantly below their full text settings. It is interesting to note that Longformer base yields much better results in the hardest setting (1 Hoax vs 100 Real) when exposed only to definitions. This is indeed a surprising and counter intuitive result that deserves future investigation.

5.1 Ablation: Effect of Definitions on Classifying Hoaxes

We run a data ablation test on the full-text split of the dataset to find the impact of definition sentences. To this end, we remove the first sentence of each article, and replicate the “full text” classification experiment, focusing on RoBERTa-Large, the most consistent model. The results of this ablation experiment are shown in Table 2, suggesting that F1 goes down about 2% for the positive class when the definition sentence is missing. This shows that definitions show critical information about entities and events in Wikipedia, but often are not the place where hoax features would emerge, and therefore removing them from the full text doesn’t change much the story.

6 Conclusion and Future Work

We have introduced HOAXPEDIA, a dataset containing hoax articles extracted from Wikipedia, from a number of sources, from official lists of hoaxes, existing datasets and Web Archive. We paired these hoax articles with similar real articles, and after analyzing their main properties (conclud-

Model	Model Size	Definition			Fulltext		
		1H2R	1H10R	1H100R	1H2R	1H10R	1H100R
Albert-base-v2	12M	0.23	0.17	0.06	0.67	0.47	0.11
Albert-large-v2	18M	0.28	0.30	0.15	0.72	0.63	0.30
BERT-base	110M	0.42	0.30	0.14	0.55	0.57	0.32
RoBERTa Base	123M	0.57	0.59	0.53	0.82	0.75	0.63
Longformer-base	149M	0.43	0.35	0.54	0.80	0.78	0.67
T5-Base	220M	0.48	0.25	0.14	0.51	0.27	0.23
BERT-large	340M	0.43	0.36	0.17	0.61	0.64	0.33
RoBERTa-large	354M	0.58	0.63	0.62	0.84	0.81	0.79
T5-large	770M	0.54	0.32	0.13	0.54	0.43	0.37

Table 1: F1 on the positive class - *hoax* at different degrees of data imbalance for Definition and Fulltext Settings(H: Hoax, R: Real)

Model	Setting	Precision	Recall	F1
RoBERTaL	Ft 1:2	0.83	0.80	0.82
RoBERTaL	Ft 1:10	0.82	0.71	0.76
RoBERTaL	Ft 1:100	0.67	0.51	0.58

Table 2: Performance of RoBERTa-Large on Binary Classifications Without Definition Sentences in Articles (with Fulltext Hoax:Real ratio in Settings column)

ing they are written with very similar style and content), we report the results of a number of binary classification experiments, where we explore the impact of (1) positive to negative ratio; and (2) going from the whole article to only the definition. This is different from previous works in that we have exclusively looked at the content of these hoax articles, rather than metadata such as traffic or longevity. For the future, we would like to further refine what are the criteria used by Wikipedia editors to detect hoax articles, and turn those insights into a ML model, and explore other types of non-obvious online vandalism.

7 Ethics Statement

This paper is in the area of online vandalism and disinformation detection, hence a sensitive topic. All data and code will be made publicly available to contribute to the advancement of the field. However, we acknowledge that deceitful content can be also used with malicious intents, and we will make it clear in any associated documentation that any dataset or model released as a result of this paper should be used for ensuring a more transparent and trustworthy Internet.

References

B. Thomas Adler, Luca de Alfaro, Santiago M. Mola-Velasco, Paolo Rosso, and Andrew G. West. 2011.

Wikipedia vandalism detection: Combining natural language, metadata, and reputation features. In *Computational Linguistics and Intelligent Text Processing*, pages 277–288, Berlin, Heidelberg. Springer Berlin Heidelberg.

Sumit Asthana and Aaron Halfaker. 2018. [With few eyes, all hoaxes are deep](#). *Proc. ACM Hum.-Comput. Interact.*, 2(CSCW).

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Quang Vinh Dang and Claudia-Lavinia Ignat. 2016. [Quality assessment of wikipedia articles without feature engineering](#). In *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*, JCDL '16, page 27–30, New York, NY, USA. Association for Computing Machinery.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings 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 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Anis Elebiary and Giovanni Luca Ciampaglia. 2023. The role of online attention in the supply of disinformation in wikipedia. *arXiv preprint arXiv:2302.08576*.

Rudolf Flesch. 2007. Flesch-kincaid readability test. *Retrieved October*, 26(3):2007.

Jim Giles. 2005. Special report internet encyclopaedias go head to head. *nature*, 438(15):900–901.

Eduard Hovy, Roberto Navigli, and Simone Paolo Ponzetto. 2013. Collaboratively built semi-structured content and artificial intelligence: The story so far. *Artificial Intelligence*, 194:2–27.

M. Hu, E. Lim, A. Sun, H. W. Lauw, and B. Vuong. 2007. [Measuring article quality in wikipedia](#). *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*.

Srijan Kumar, Robert West, and Jure Leskovec. 2016. Disinformation on the web: Impact, characteristics, and detection of wikipedia hoaxes. In *Proceedings of the 25th International World Wide Web Conference*.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

Jodi Schneider, Bluma S. Gelley, and Aaron Halfaker. 2014. [Accept, decline, postpone: How newcomer productivity is reduced in english wikipedia by pre-publication review](#). In *Proceedings of The International Symposium on Open Collaboration, OpenSym '14*, page 1–10, New York, NY, USA. Association for Computing Machinery.

Arsim Susuri, Mentor Hamiti, and Agni Dika. 2017. [Detection of vandalism in wikipedia using metadata features – implementation in simple english and albanian sections](#). *Advances in Science, Technology and Engineering Systems Journal*, 2:1–7.

William Yang Wang and Kathleen McKeown. 2010. [“got you!”: Automatic vandalism detection in Wikipedia with web-based shallow syntactic-semantic modeling](#). In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 1146–1154, Beijing, China. Coling 2010 Organizing Committee.

KayYen Wong, Miriam Redi, and Diego Saez-Trumper. 2021. [Wiki-reliability: A large scale dataset for content reliability on wikipedia](#). In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21*, page 2437–2442, New York, NY, USA. Association for Computing Machinery.

Honglei Zeng, Maher A. Alhossaini, Li Ding, Richard Fikes, and Deborah L. McGuinness. 2006. [Computing trust from revision history](#). In *Proceedings of the 2006 International Conference on Privacy, Security and Trust: Bridge the Gap Between PST Technologies and Business Services, PST '06*, New York, NY, USA. Association for Computing Machinery.

Dataset Setting	Dataset Type	Split	Number of Instances		
			Non-hoax	Hoax	Total
1Hoax2Real	Definition	Train	426	206	632
		Test	179	93	272
1Hoax2Real	Full Text	Train	456	232	688
		Test	200	96	296
1Hoax10Real	Definition	Train	2225	203	2428
		Test	940	104	1044
1Hoax10Real	Full Text	Train	2306	218	2524
		Test	973	110	1083
1Hoax100Real	Definition	Train	20419	217	20636
		Test	8761	82	8843
1Hoax100Real	Full Text	Train	22274	222	22496
		Test	9534	106	9640

Table 3: Dataset details in different settings and splits

A Appendix: Dataset Details

We release our dataset in 3 settings as mentioned in Section 4. The settings with data splits and their corresponding sizes are mentioned in Table 3.

B Appendix: Training Details

We train our models with the configuration given below. We use one NVIDIA RTX4090, one NVIDIA V100 and one NVIDIA A100 GPU to train the models.

- Learning rate: 2e-06
- Batch size: 4 (for Fulltext experiments) and 8 (For Definition experiments)
- Epochs: 30
- Loss Function: Weighted Cross Entropy Loss
- Gradient Accumulation Steps: 4
- Warm-up steps: 100