

Machine Learning Project Report: Fake News Detection

Philipp Hoffmann, Martin Wagner
Applied ML – Ludwig-Maximilians-Universität München
philipp.hoffmann@campus.lmu.de, wagner.mar@campus.lmu.de

July 16, 2025

1 Task Overview

Due to the enormous amount of news published every day, distinguishing between fake and real news is becoming an increasingly difficult task for humans, making the use of computer models indispensable. In this project, we tested different models and text preprocessing methods using the "fake-and-real-news-dataset"¹. This dataset collects titles and texts of more than 40000 news articles classified as fake or true news.

2 Methods

We compared the performance of logistic regression and linear SVM models (mainly on the collection of titles) using different vectorization techniques from natural language processing. We used `TfidfVectorizer` and `CountVectorizer` from `scikit-learn` for basic text vectorization, which also handle tokenization, including lower-casing and punctuation removal.

Given a vocabulary of tokens (e.g. words), bag-of-words assigns to a document d for each token t , the term-frequency $\text{tf}(t, d)$, defined as the count of term t in d divided by the total number of terms.² A bit more evolved is TF-IDF vectorization, which also takes into account the frequency of how often a term occurs in a fixed corpus $D = \{d_1, \dots, d_n\}$ of documents d_1, \dots, d_n . Given a term t , the inverse-document frequency of t in D is $\text{idf}(t, D) := 1 + \log\left(\frac{1+n}{1+\text{df}(t)}\right)$, where $\text{df}(t)$ is the number of documents in D containing the term t .³ Then the TF-IDF-score of a term t with respect to the collection D is defined as

$$\text{TF-IDF}(t, d, D) := \text{tf}(t, d) \text{idf}(t, D).$$

Note that d does not need to be contained in the collection D of documents itself. We applied $L2$ -normalization to the obtained vectors and additionally applied Z-score normalization to the resulting feature matrices, as this yielded better results, particularly for linear SVM. We adopted models in `couselib` to support sparse matrix operations and implemented methods to adapt to feature shifts (due to Z-score normalization) without destroying the computational efficiency resulting from the sparsity of the unshifted feature matrices. Moreover, we implemented custom tokenizers using `nltk` that support stemming and lemmatization, i.e. reduce words to a more basic form. To test all these options more compactly, we implemented a multi-column vectorizer that supports vectorization of text data from one or multiple columns of a dataframe with various vectorization options.

3 Experiments and Results

As both bag-of-words and TF-IDF yield highly sparse, high-dimensional vectors, one expects the problem to be linearly separable, thus we restricted our experiments to linear models, namely logistic regression and linear SVM. For all our experiments, we used an 80/20 train-test split. We labeled true news with 1 and fake news with 0 (for logistic regression) or -1 (for linear SVM), respectively. All of our models were trained for a total of 100 epochs.

We compared the models on TF-IDF vectorized text data with good-performing learning rates with full batch size. While both models achieved almost perfect train accuracy, logistic regression performed slightly better and faster. Smaller batch sizes could not improve results, but slowed down the training time significantly (see Figure 1). Because of this, we only used the logistic regression model with full batch size for further experimentation.

¹<https://www.kaggle.com/datasets/clmentbisaillon/fake-and-real-news-dataset>

²Usually bag-of-words is defined without dividing by the total number of terms, but this does not matter for us, as we normalize anyway

³Often times one also defines $\text{idf}(t, D) := \log\left(\frac{n}{1+\text{df}(t)}\right)$, but we stick with the convention implemented in `scikit-learn`

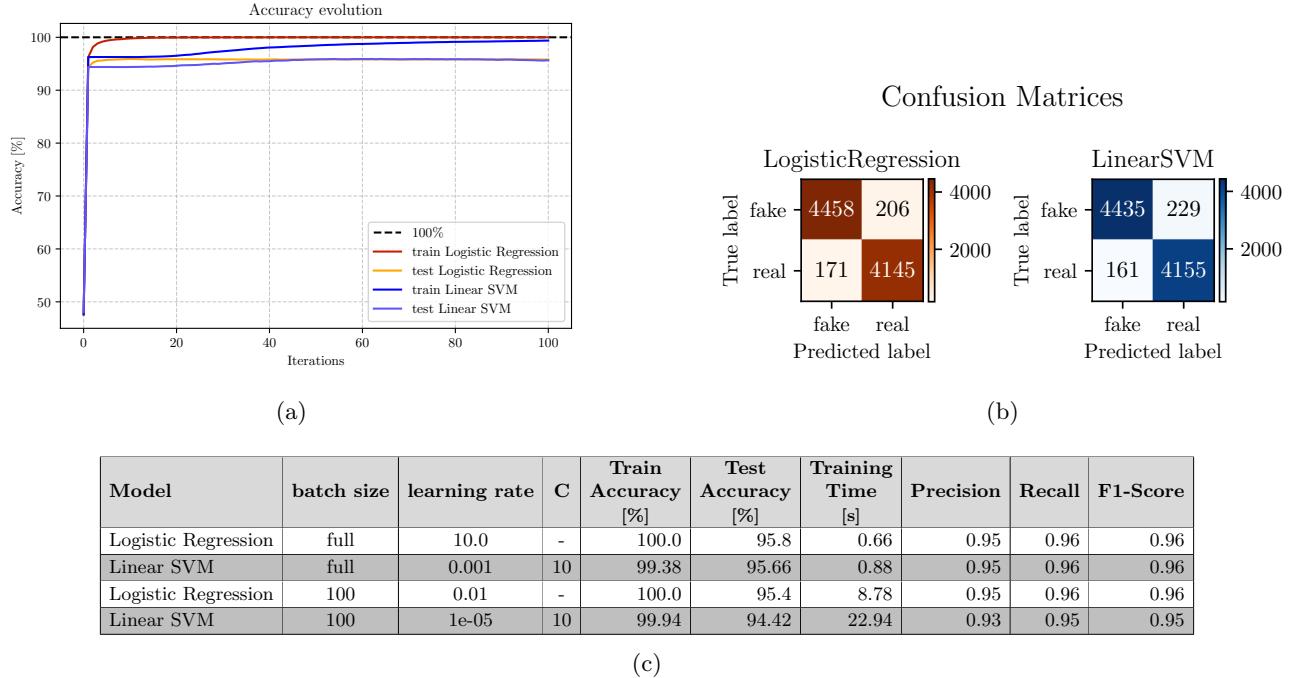


Figure 1: Model comparison between logistic regression and linear SVM. (a): Accuracy evolution on full batch, (b): Corresponding confusion matrices, (c): Comparison of mini batch and full batch

For the second part of our experiment, we tested different vectorization methods. Figure 2(a) shows how the number of features affects the accuracy: In general, more features result in better accuracy, but more than 5000 features yield only a minor improvement. This trend is also confirmed by Figure 2(b): When using text or text and titles as data, the number of features increases significantly, while also resulting in even better test accuracy.

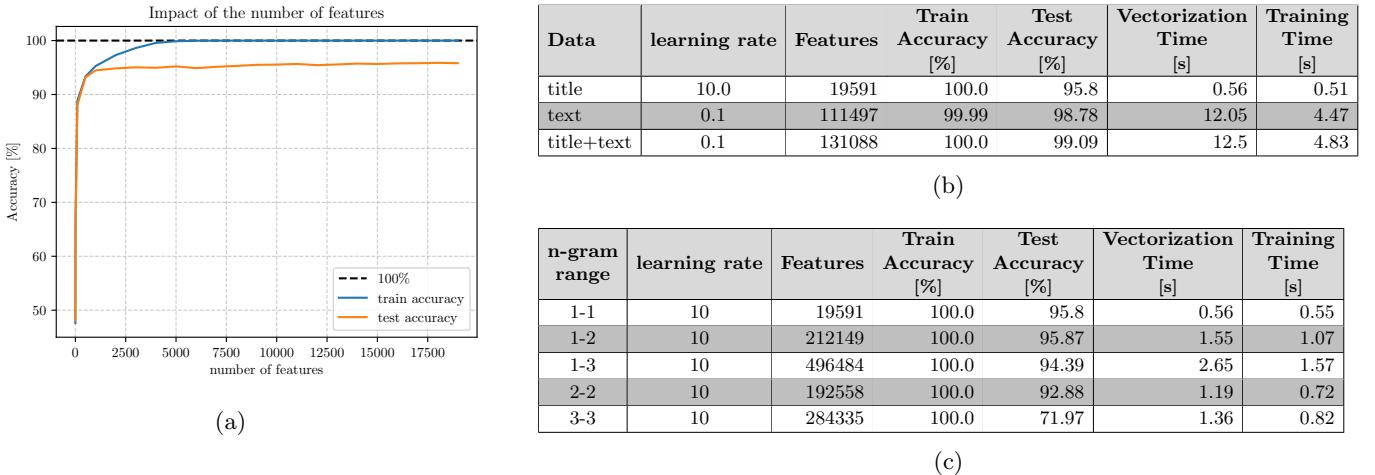


Figure 2: (a): Impact of number of features, (b): Performance with different data, (c): Influence of n-gram ranges, n-gram range given as minimum length - maximum length

With that in mind, we compared TF-IDF as well as bag-of-words vectorization using four different tokenization methods, with and without stop word removal. As Table 1 highlights, all combinations achieved similar performance, but stop word removal slightly decreases test accuracy. Moreover, our custom tokenizers improved the test accuracy slightly, but the vectorization time increased heavily with the complexity of the tokenizer.

Vectorization	Stop Words	Tokenizer	Features	Train Accuracy [%]	Test Accuracy [%]	Precision	Recall	F1-Score	Vectorization Time [s]
tf-idf	None	default	19591	100.0	95.8	0.95	0.96	0.96	0.59
tf-idf	None	basic	24931	100.0	97.18	0.97	0.98	0.97	5.07
tf-idf	None	stemming	17387	100.0	96.78	0.96	0.97	0.97	11.77
tf-idf	None	lemmatization	20916	100.0	96.84	0.96	0.97	0.97	33.35
tf-idf	english	default	19319	100.0	94.8	0.94	0.95	0.95	0.57
tf-idf	english	basic	24661	100.0	96.36	0.96	0.97	0.96	5.07
tf-idf	english	stemming	17247	99.99	96.28	0.96	0.96	0.96	10.66
tf-idf	english	lemmatization	20871	100.0	96.16	0.96	0.96	0.96	30.73
bag-of-words	None	default	19591	100.0	95.8	0.95	0.96	0.96	0.57
bag-of-words	None	basic	24931	100.0	97.12	0.97	0.97	0.97	5.05
bag-of-words	None	stemming	17387	100.0	96.84	0.96	0.97	0.97	11.63
bag-of-words	None	lemmatization	20916	100.0	96.9	0.97	0.97	0.97	33.67
bag-of-words	english	default	19319	99.99	94.8	0.94	0.95	0.95	0.59
bag-of-words	english	basic	24661	100.0	96.48	0.96	0.97	0.96	5.06
bag-of-words	english	stemming	17247	99.99	96.38	0.96	0.96	0.96	10.61
bag-of-words	english	lemmatization	20871	100.0	96.21	0.96	0.96	0.96	29.63

Table 1: Comparison of different vectorization configurations with four different tokenizers. default: default tokenization from `TfidfVectorizer/CountVectorizer` without custom tokenizer, basic: `nltk.word_tokenize` combined with punctuation removal and lowercasing, stemming: basic tokenizer combined with `SnowballStemmer` from `nltk`, lemmatization: basic tokenizer combined with `WordNetLemmatizer` from `nltk`. (used learning rate: 10)

Lastly, we investigated which words or short phrases the model identified as most relevant for distinguishing between real and fake news (Figure 3). For this, we trained the model on different combinations of n-grams. Figure 2(c) shows that training only with bi- and especially trigrams significantly decreased the test accuracy.

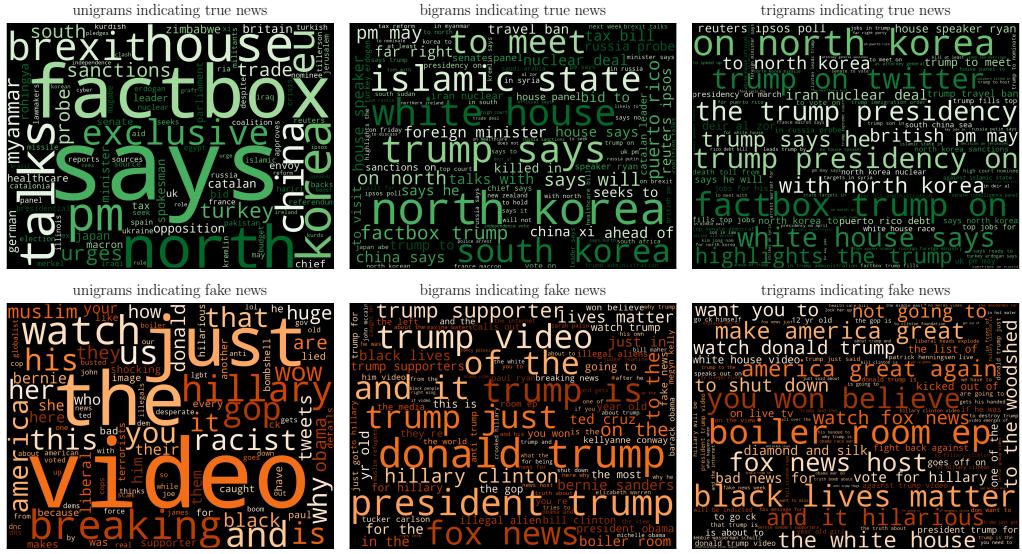


Figure 3: Illustration of the most important n-grams identified by the model trained only on uni-, bi-, trigrams resp.

4 Discussion

Our experiment shows that linear models are indeed well-suited for detecting fake news, even when analyzing only news titles. If one also trains the model on the whole body of text, almost perfect classification is achievable. Removing stop words from the titles seems to worsen the performance, which could indicate that titles become too short or that they are indeed more common in either fake or real news. It seems like machine learning models prefer single words for classification, whereas humans usually need context for understanding text. However, our models only give binary classification and cannot detect the exact passages containing fake news, which could be an interesting objective for future research. A starting idea for this could be to divide the text into smaller paragraphs.