

Deepfake tweets detection

Project Proposal for NLP Course, Winter 2023

Adam Frej

Warsaw University of Technology
01151392@pw.edu.pl

Adrian Kamiński

Warsaw University of Technology
01151387@pw.edu.pl

Piotr Marciniak

Warsaw University of Technology
01151428@pw.edu.pl

Szymon Szmajdziński

Warsaw University of Technology
01151438@pw.edu.pl

supervisor: Anna Wróblewska

Warsaw University of Technology
anna.wroblewska1@pw.edu.pl

Abstract

This research project aims to address the growing concern of deepfake text generation, specifically focusing on detecting deepfake tweets. Deepfakes, generated using advanced machine learning techniques, have the potential to spread misinformation, impersonate individuals, and facilitate malicious activities. Detecting such deepfake content is crucial to protect society from deception and harm. The project leverages the TweepFake dataset, which contains tweets from humans and bots. It also explores various text representations and preprocessing techniques.

The research questions revolve around the development of a reliable deepfake detection algorithm. The project hypotheses involve exploring patterns, such as the use of emoticons, mentions, and misspelt words, that may indicate machine-generated tweets. Different machine learning and deep learning models were used to maximize detection accuracy while maintaining precision and recall balance. Our top-performing model achieved an 87% accuracy rate in detecting deepfakes across all categories and an 83% accuracy rate specifically in identifying deepfakes generated by

GPT-2.

The proposed work intends to improve state-of-the-art deepfake tweet detection and by creating well performing model. We believe that our detection model will contribute to more trustworthy online interactions by detecting deepfakes on the internet. Ultimately, the project seeks to enhance users' safety, trust, and confidence in their online experiences.

1 Introduction

The goal of our project is to determine if the content of a tweet is a deepfake. The term “deepfake” is a portmanteau of “deep learning” and “fakes” (Vincent, 2018). It refers to a type of synthetic content that is created using advanced machine learning techniques, particularly deep learning algorithms. Deepfakes are typically associated with manipulated videos but can also involve audio, images, and text. The core characteristic of a deepfake is that it convincingly alters or generates content to make it appear as if it is authentic, even when it is not. This work focuses on deepfake related to text corpora, in particular, to tweet data. This depicts human interaction in a short, dynamic form. ~~Texts~~ Tweet texts tend to be shorter and contain less context, which poses more challenges.

There are several reasons to tackle the problem of deepfake tweet detection. One of them

can be protection against misinformation. Deepfake tweets can be used to spread misinformation ([wrong information](#)), [disinformation \(false information spread to deceive\)](#), ~~disinformation~~, and fake news ~~—(false stories that appear to be news, spread on the internet or using other media, usually created to influence political views or as a joke)~~, which can have serious real-world consequences, such as influencing public opinion ([Milmo, 2023](#)), election outcomes ([Lynn, 2024](#)), and even inciting violence ([Parsons, 2022](#))([definitions in brackets come from Cambridge Dictionary](#)). Detecting and mitigating deepfakes is crucial to protect society from being misled by false narratives. Moreover, an increasing number of synthetic tweets can lead to a lack of trust in information shared on social media platforms.

Another reason for tackling the problem of deepfake tweet detection is the protection of individuals against impersonation ([Fernandez, 2024; Parsons, 2022](#)), privacy violations ([Rahman-Jones, 2024](#)), and identity theft ([Europol Innovation Lab, 2024](#)). Deepfake tweets can be used to impersonate both public figures and private citizens, and by addressing this issue, we can safeguard people's rights and personal information.

Deepfakes can also be exploited for malicious purposes, including extortion ([Fernandez, 2024](#)), [fraud \(Europol Innovation Lab, 2024\)](#), ~~fraud~~, and harassment ([Rahman-Jones, 2024](#)). They can be used to create a network of fake users who encourage other users to use some services ([Fernandez, 2024](#)). Detecting deepfakes is necessary to prevent these harmful actions.

Detection of deepfake corpora can also lead to exposure of flaws in existing text-generating algorithms, which can help to improve the language models in the future. This means creating texts that are more exciting and semantically plausible to read. Furthermore, comparing the models to human benchmarks can identify the types of errors in generated sequences that humans tend to notice and make text appear unnatural.

Properly detecting deepfake tweets can increase users' trust, safety, and confidence in online interactions and content consumption. It can contribute to a more positive online experience.

2 Related works

The TweepFake dataset was created in (Fagni et al., 2021), where authors compared multiple approaches in both corpora encoding and algorithms to detect fake tweets. What is unique about this dataset is that the data used there was generated by various algorithms and was extracted directly from social media. The following setup can provide more generic results when evaluating models. Their ~~results suggest research suggests~~ that transformer-based models offer better results. [The best result was achieved by Roberta \(Liu et al., 2019\) \(89.6% accuracy\)](#). They also found out that all [their](#) approaches struggled with tweets generated by GPT-2 (Radford et al., 2019), and the best results ~~were obtained using on~~ [GPT-2 tweets were obtained with](#) an RNN decoder using character encoding ([82% of accuracy vs 74% accuracy achieved by the best model overall](#)).

One of the main difficulties in natural language processing is text representation. Related works investigate several solutions to this problem in the context of feeding data models detecting machine-generated corpora. The TweepFake (Fagni et al., 2021) article proposes 4 options. One of them is a popular method, bag-of-words (BoW), with features weighted using TF-IDF function (Sebastiani, 2002). The output of this methodology was processed by either logistic regression (Cox, 1958), random forest (Ho, 1995), or SVM (Cortes and Vapnik, 1995). However, this approach suffers from the curse of dimensionality, as the features are very sparse and require a lot of data. Another drawback is that the algorithm misses the semantic context of the words. The article got the worst result here, hovering around 0.80 accuracy.

Another option is to encode text using a high-level language model, which overcomes these limitations. The TweepFake [article](#) used BERT (Devlin et al., 2019) to provide contextual embeddings that include words context and can be encoded into a vector representing specific text. Yet again, this representation is later processed by classifiers.

The third approach operates on the character level. A vocabulary of characters is mapped to the internal embedding matrix, which is passed to the selected deep-learning networks. This results in a surprisingly good score, up to 0.85 accuracy, and can be useful in cases without ~~pretrained~~ [pre-trained](#) models available, for example, another language.

However, the most successful approach is utilizing pre-trained models. They take raw input directly and, with fine-tuning on a given dataset, solve the classification problem of labeling labelling a sentence as human or machine-generated. This way, the models operate on the complex sequence-based understanding level. The TweepFake article tested several language models, all related to BERT. Besides the original BERT, XLNet (Yang et al., 2020) and RoBERTa (Liu et al., 2019) were utilized as they can produce 15% better results thanks to architecture modifications and a bigger training dataset. The article also tested DistilBERT (Sanh et al., 2020), which tries to keep performance while simplifying the architecture and halving the number of parameters. Generally, ~~the results~~ most results of the tested methods get up to 0.90 accuracy, the best being RoBERTa, which indicates that the most complex representations are most effective. They are also well-balanced in terms of precision and recall.

One of the works tried to detect generated text by exploiting situations when humans are fooled by it (Ippolito et al., 2020). They also investigated several text representations. Similarly, the primary one is a fine-tuned BERT (Devlin et al., 2019). Again, this approach far surpassed other methods depicted in the article.

Another representation in this work is a simple BoW. This time, the GPT-2's 50,000 token vocabulary (Sennrich et al., 2016) is used to count the occurrences of tokens in text sequences. This embedding allowed for training logistic regression binary classifier, which achieved the next best result. ~~but the differences in outcomes were huge.~~ Depending on the dataset, we could observe 9, 20, and 21 percentage point differences in accuracy, where a 9 pp difference was on the easiest dataset (BERT achieved the best result on it).

The article also proposed Histogram-of-Likelihood Ranks. ~~As in GLTR~~ Similarly to another detection tool - GLTR introduced in (Gehrmann et al., 2019), they created an energy-based deepfake text decoder by calculating the probability distribution of the next word given the previous words in the sequence according to the GPT-2 language model. They ranked the words by likelihood and then binned them either into 4 groups or uniformly over the whole vocabulary. Such histograms served as input for logistic regression binary classifiers. Un-

fortunately, this approach resulted in worse scores despite trying to reproduce successes reported in GLTR, ~~which can~~. The authors suggest it may be explained by training data selection.

Another deepfake detection research performed on social media texts was conducted on Amazon reviews in article (Adelani et al., 2019). Fake news ~~were~~ was generated using the GPT-2 text generation model (Radford et al., 2019). Authors, in order to adjust the model to new corpora, adapted the original GPT-2 model to Amazon (He and McAuley, 2016) and Yelp reviews (Zhang et al., 2015). As for detection algorithms, they used Grover (~~Zellers et al., 2019a~~), ~~GLTR~~ (Zellers et al., 2019a) (neural network trained in a semi-supervised way to detect if the content was generated by bot), GLTR (Gehrmann et al., 2019), and OpenAI GPT-2 (Solaiman et al., 2019). They also tried combining those models using logistic regression at the score level. One of the experiments they performed was selecting a real review out of 4, of which 3 were fake. Human participants tended to randomly guess which review was real since they were right about 25% of the time. However, the models were not much better. The best configuration achieved 20% accuracy on this task.

In (Guo et al., 2023), authors proposed an HC3 dataset consisting of questions and their corresponding human/ChatGPT answers. Based on this dataset, they conducted a comprehensive human evaluation and linguistic analysis, as well as developed several detecting models. They used the GLTR model with logistic regression, RoBERTa, and RoBERTa-QA - a Question Answering version of the model that supports a text pair input format, where a separating token is used to join a question and its corresponding answer. In their work, they came to the following conclusions:

- The robustness of the RoBERTa-based-detector is better than GLTR.
- RoBERTa is not affected by indicating words (characteristic words for ChatGPT).
- RoBERTa is effective in handling Out-Of-Distribution scenarios, ~~whereas we can observe a significant decrease in performance on GLTR's when testing on data in first-seen format.~~
- Detecting ChatGPT-generated texts is more difficult in a single sentence than in a full text.

3 Datasets

~~Several datasets were created to build a solution to detect content created by deep learning algorithms. Some of them are presented below. Below, we present the datasets we used for our research. We extended the TweepFake data with GPT-2 related texts.~~

1. The TweepFake dataset (Fagni et al., 2021) – it is a dataset which contains 25,572 tweets half human and half bots generated. They ~~used~~ selected 17 human accounts, which were imitated by the 23 bots. The selection of bot accounts was based on Twitter’s profile descriptions or profile URLs or related GitHub. Some of the fake accounts imitated the same human profile. These bots were using different technologies, and for almost all of them ~~(except one bot account), the used technology is known, the technology they used is known (except for two bot accounts).~~
2. The GPT-2 output datasets (Ippolito et al., 2020) ~~— it is a group of several datasets in which deepfakes are generated by GPT-2 models (Radford et al., 2019). The datasets differ because different decoding strategy settings and different sizes of models are applied during generation. Each dataset contains 500,000 training and 5,000 validation and test samples, which are evenly spread across classes (human-generated excerpts of web texts or GPT-2 generated).~~
3. ~~The HC3-English datasets (Guo et al., 2023) — it is a group of datasets from different sources, in which for each question, there is provided at least one human, and ChatGPT3.5 answer. The questions and answers by the human experts come mainly from publicly available question-answering datasets. There is also an additional source in which Wikipedia is treated as a human expert who is asked questions based on concepts in crawled data.~~

In most papers (Zellers et al., 2019b; Bakhtin et al., 2019), datasets containing generated and human texts are not provided because big corpora are used to build a generative model ~~in which~~ descriptors where discriminators act as detectors of deepfakes. Later, they test their ~~descriptors~~ discriminator on unused parts of corpora and the

text produced by ~~generators~~ a generator to see if the ~~descriptors work~~ discriminator works correctly.

4 Approach & research methodology

Our research on detecting deep fakes of tweets was performed mainly on the TweepFake dataset. We also ~~tried to use one of~~ introduced the GPT-2 output datasets to improve the detection of deep fakes prepared by the GPT-2, ~~which.~~ This caused the most problems in the original work, which we investigated further. To test our solution, we used the split provided by the authors of the dataset. They performed a stratified split into 3 datasets (training, validation, test), which ~~roughly reminds~~ approximates the 80%-10%-10% ~~By doing it this way, later,~~ division (precisely numbers of samples in each subset: 20,712 - 2,302 - 2,558). They also defined 3 categories of generative content based on two main methods (RNN, GPT-2) and aggregated the rest to others category. By having the same settings, we can compare our results with the ones achieved in the article.

4.1 Initial research

In our project, we stated several research questions, which set the direction of the experiments:

- Can we build a reliable deepfake detection algorithm? By reliable algorithm, we mean the model that maximises accuracy on a balanced dataset while remaining well-balanced in precision and recall. That means detecting generated tweets while avoiding assigning false positives.
- What are the most effective embeddings for deepfake detection in tweets?
- Are there any patterns that indicate the model-generated tweet content?

There are also some hypotheses which we tested:

- Emoticons: The use of emoticons may be higher in human-generated content.
- Mentions: The use of mentions of other users may be higher in human-generated content.
- Misspells: There will be more misspelt words in content generated by bots.

- URLs: The impact of different URL encoding, e.g., encoding all URLs to a single token vs extracting the basepath of the URLs.

In order to investigate and answer these hypotheses, we conducted explanatory data analysis (EDA). This also allowed us to gain insight into data and conduct further preprocessing steps. We assessed the most basic data properties. ~~All~~ The original splits are balanced, both in terms of human ~~→~~ non-human content and different generative methods in a non-human class. The tweet lengths are also similar across different classes. Using the Natural Language Toolkit (NLTK) (Bird et al., 2009), we tokenized the tweets and counted word occurrences. As expected, the most frequent are stop words (e.g. *the, to, a*), which we removed for later research. Lastly, the sentiment analysis of tweets reports ~~is that they are~~ quite balanced and ~~focused around neutral values of the polarity and subjectivity accumulated around neutral polarity values (around 0)~~ in different classes ~~with regard to splits and splits, but human-created content tends to be more subjective than the one created by bots (median of subjectivity is higher in case of human content)~~. For this, we used *SpacyTextBlob* ~~from extension to the~~ spaCy package (Honnibal et al., 2020).

Regarding the hypotheses, we checked them by ~~simple~~ aggregations performed on texts. ~~The results are available in the appendix and are referenced below.~~

1. **Emoticons** - Figure 3: most tweet do not contain any. However, if some appeared, they were mostly human-generated, ~~which roughly confirms the hypothesis~~. Importantly, the main generative methods ~~gpt2 and rnn~~ gpt2 and rnn, were least likely to use any. There were some algorithms in ~~other~~ others category which did use emojis.
2. **Mentions** - Figure 4: again, most tweets do not contain them. Here, the difference is more apparent. Humans are almost the only ones to use mentions, directly confirming the hypothesis.
3. **Misspells** - Figure 5: they are mainly introduced by humans. However, the situation is more balanced. The tweets without mistakes are not dominant and are mainly made by bots. Generally, the more mistakes in a

tweet, the more likely it is to be created by a human. The hypothesis appears to be untrue. We used *pyspellchecker* ~~(Barrus, 2019) with basic English dictionary en-core-web-sm from spaCy package (Honnibal et al., 2020)~~ (Barrus, 2019). Perhaps the human texts contain more slang and informal speech, which technically contains misspells.

4. **URLs** ~~÷ similarly~~ Figure 6: similar to mentions, most tweets do not contain them, and humans generate most URLs. Only ~~rnn~~ RNN tried to introduce some. This makes the hypothesis obsolete, as the impact of encoding does not matter when generated texts can not be encoded.

Generally, excluding misspells, all the discussed properties rarely appear. Usually, the tweets do not contain them. Therefore, despite assessing the hypotheses as true or not, they are not very significant to our research.

~~We applied some simple preprocessing steps to~~

4.2 Preprocessing

~~To the texts of the tweets. Using again NLTK and other basic Python libraries, we applied preprocessing steps, which are presented below. Using the NLTK tokenizer and REGEXes again, we tagged URLs and mentions. As reported in EDA, they did not contain meaningful information. (replaced URL and mentions with the corresponding words <URL> and <MENTION>).~~ Then, we tokenized the tweets and removed stop words. We also created two versions of the data, with applied stemming or lemmatization. Both are later tested.

As a result of this preprocessing, we ended up with 4 versions of the Tweep Fake dataset:

- raw – Raw version of the dataset without any preprocessing (this dataset was utilized only on transformer models),
- simple – raw dataset after performing mentions and URL tagging and removing the stop words,
- lemmatized – simple dataset after performing lemmatization,
- stemmed – simple dataset after performing stemming,

- [bert_embeddings – embeddings extracted from BERT on simple dataset.](#)

We explored ~~five different~~ three different groups of approaches to finding the best model for this task~~—~~, which are more precisely described in the corresponding subsections:

- [Machine learning approach \(subsection 4.3\) – In this approach, we tested five classifiers with different text encoding methods.](#)
- [Deep learning approach \(subsection 4.4\) – We tested two approaches for tokenization of tweets. One was on the level of characters, and the second was on the level of words. For each level of tokenization, we tested three different architectures.](#)
- [Transformers \(subsection 4.5\) – We tested three different transformers architectures with different settings.](#)

4.3 [Machine learning models](#)

The first approach is based on text representation with a bag of words encoded with the TF-IDF function. Next, tweets encoded this way were processed by machine-learning algorithms. In this project, we use five well-known classifiers: Light Gradient Boosting Machine (LightGBM) (Ke et al., 2017), eXtreme Gradient Boosting (XGBoost) (Chen and Guestrin, 2016), Random Forest (Ho, 1995), Logistic Regression (Cox, 1958), Support Vectors Machine (SVM) (Cortes and Vapnik, 1995).

As the first approach often suffers from the curse of dimensionality because feature space after encoding is sparse and does not take into account the word order, in the second approach, we prepared the feature space using BERT (Devlin et al., 2019). BERT provides contextual embeddings, i.e. fixed-size vectors representing words which depend on the context in which the word occurs. We averaged these embeddings to obtain a fixed-size vector of a tweet text, which depends on context and word order [and can be used as the embedding of the tweet.](#) As in the previous scenario, the tweets encoded by BERT have been later used to learn the same set of classifiers.

In both the first two approaches, the best hyperparameters for our ML pipeline (we optimized parameters also for the TF-IDF Vectorizer) were found with the help of the Optuna package (Akiba

et al., 2019). As we wanted to get the best hyperparameters for our problem in these approaches, we merged our training and validation dataset to obtain a more extensive dataset. Later, during the search for hyperparameters, this bigger dataset was split into five stratified folds. Optuna was trying to optimize the mean of the balanced accuracy metric calculated on the test fold, which, with the stratified approach, should have very close results to simple accuracy. We set some restrictions within Optuna. Each ML pipeline takes 20 minutes, and maximally, 100 tries to find the optimal hyperparameters. The hyperparameters and their value space are presented in Table 4.

4.4 [Deep learning approach](#)

In the third and fourth approaches, we used deep learning techniques to build the models to solve our problem. In the third scenario, we encoded text by working at the character level. We mapped the text of the tweets to lowercase and tokenized the letters. After tokenization, we padded sequences to the maximum length of 320 characters ~~—because even though the maximal tweet length right now is 280 characters in the original dataset, they were longer tweets. One reason may be that the tweets used contained HTML tags that are not counted as characters in posted tweets.~~

In the fourth scenario, we applied token encoding of the text. We mapped the text of the tweets to lowercase and removed punctuation. The tokenizer was using the ~~15,000~~ 15,000, the most occurring tokens and the maximal length of the sequence was 100 tokens.

In both DL approaches, we used similar architectures of deep neural networks. They differed only in the number of filters in the convolution layer and the number of units in the GRU layer. The approach with tokens, in most cases, used the double of filters/units in their architectures. We checked three types of architectures:

- CNN ~~—in—~~ [In](#) this architecture, convolutional layers on the activation matrix coming from the embedding layer are used; after this layer, we performed global max pooling and used a dense layer to classify,
- GRU ~~—in—~~ [In](#) this architecture, we used only a bidirectional GRU layer on the activation matrix; after this layer, we performed global max pooling and used a dense layer to classify,

- CNN+GRU ~~→ in~~ In this architecture, we combined the strength of both architectures; we concatenated the outputs from the CNN layers and GRU layer to perform the classification.

In all DL approaches, we shared the following configuration:

- batch size: 256
- epochs: 200 with early-stopping (10 epochs patience)
- save: ~~model from best (in regard to loss function) and last evaluation step~~ the best model (regarding loss function on validation dataset)
- optimizer: Adam
- learning rate: 1e-4
- loss: CrossEntropy

Due to computation limitations and the high demand for resources of deep learning methods, we did not perform hyperparameter optimization with these kinds of models.

4.5 Transformers

In the last approach, we focused on larger models. Finally, we tried the following architectures:

- xlm-roberta-base (Conneau et al., 2019) - XLM-RoBERTa-Base is a cross-lingual pre-trained language model based on the RoBERTa architecture, designed to understand and generate text in multiple languages with improved performance and efficiency.
- distilbert-base-uncased (Sanh et al., 2020) - DistilBERT-Base-Uncased is a distilled version of BERT, optimized for speed and efficiency while retaining much of the original model's performance.
- gpt2 (Radford et al., 2019) - GPT-2, or Generative Pre-trained Transformer 2, is a powerful autoregressive language model that excels in generating coherent and contextually relevant text.

The experiments shared the following configuration (unless otherwise stated in the method description):

- batch size: 8
- epochs: 10 with early-stopping (2 epochs patience)
- save: model from best (in regard to loss function) and last evaluation step
- weight decay: 0.01
- optimizer: AdamW
- loss: CrossEntropy
- seed: 1

Due to computation limitations and the high demand for resources (both compute and time) of these methods, we did not perform hyperparameter optimization with these kinds of models (other than in simple experiments described below).

We started with XLM-RoBERTa and in this approach, we trained the network 3 times, each time switching hyperparameters. First, we tried freezing layers that were not responsible for final classification (this approach will be denoted by XLM0) and without it (XLM1). After obtaining results, it was clear that freezing so many layers did not help model performance, so we discarded this idea from further training configurations. For XLM-RoBERTa, we additionally tried another value learning rate hyperparameter. In XLM1, the value $2 \cdot 10^{-5}$ was used and in new approach XLM2, we tried lower value i.e. $2 \cdot 10^{-6}$. The second approach proved to be better (by over 1.2% in accuracy), and that value was used in all later experiments.

After that, we moved to the DistilBERT model. Since this model is relatively small and fast, we decided to check how the learning process would look when we included the GPT2 output dataset to pretrain the model. In the pretraining process, we had to limit the number of steps (evaluation of the whole dataset once would take over 6 hours). We ended up with one epoch and evaluation after ~~2,000~~ 2,000 steps and early stopping with (5 evaluations step patience). However, the additional pretraining on this dataset did not improve the results, so we didn't repeat this process again on larger models.

The last architecture we trained was GPT2. We trained this model once with a setup that proved best (the one stated as a shared configuration) in previous examples.

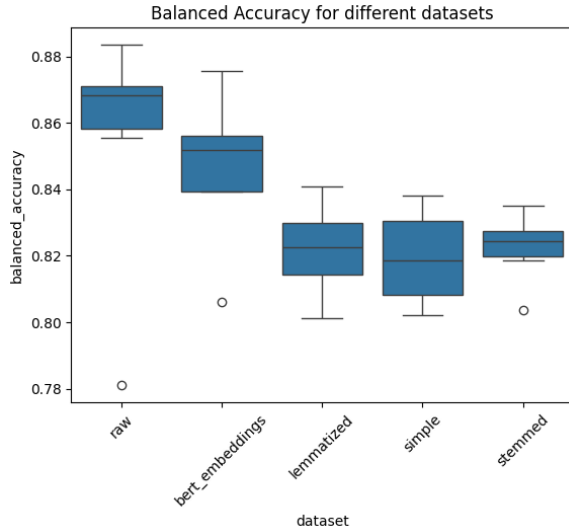


Figure 1: Balanced Accuracy for different datasets ~~embedding methods described in section 4~~

5 Results

To evaluate the results of the models described above, we used several metrics, i.e., accuracy, balanced accuracy, precision, recall and f1. The use of balanced accuracy was motivated by the fact that different labels were of varying difficulty to classify. A model that was very proficient at classifying deepfakes generated by RNNs could perform poorly in classifying deepfakes generated by GPT-2, yet still achieve reasonably high accuracy. We decided to use balanced accuracy to better represent the model’s overall performance across all classes. The evaluation was done on the testing dataset ($\sim 10\%$ of the whole dataset that was not used during the training/validation process).

Figure 1 presents the results of different models with respect to the type of data used to train them. We can see that raw data achieved the best results. This data was only used by transformers. Bert embedding has also achieved high results with respect to the accuracy score. The other 3 datasets, i.e., lemmatized, simple and stemmed, used either TF-IDF or DL embedding created only on the data from TweepFake. We can see that those embeddings proved to perform significantly worse.

Table 1 illustrates the balanced accuracy of the best 10 models. The transformers models (XLM, Distil-Bert, GPT2) achieved the most favorable outcomes when trained on the original dataset.

Notably, the second-best performing model was the SVC model utilizing embeddings generated by Bert. In a comparison with (Fagni et al., 2021), our findings align closely, with the transformer models’ results proving comparable to those reported in the referenced article. Notably, we achieved superior results with the SVC model, ~~possibly attributable to variations in data preprocessing methods.~~

Table 2 provides insights into the overall accuracy and specific accuracies categorized by the method used to create the deepfakes. Notably, detecting deepfakes generated by the GPT2 model was the most challenging task. Tweets generated by humans ranked as the second most challenging category, although their accuracies were significantly higher than those for GPT2-generated tweets. On the other hand, tweets generated using RNNs and alternative methods are generally easy to detect, having notably high accuracies.

The XML model trained with a lower learning rate (XLM2) achieved the highest global accuracy. Interestingly, the XML model without a lower learning rate (XML1) yielded the best results for identifying GPT2-generated deepfakes, while maintaining relatively high scores for other methods.

The detailed accuracies for each class type and model are presented in Figure 2. Notably, RNN and other classes emerge as the easiest to detect, with accuracies around 95% for RNNs and 92% for other class types. The heatmap further highlights the challenge in detecting deepfakes generated by GPT2, as they exhibit lower detection rates. Transformers, on the other hand, consistently showcase the highest accuracies within this class type.

Compared to the findings presented in the paper (Fagni et al., 2021), our study achieved superior results in detecting deepfakes generated by GPT2, all while maintaining a high overall accuracy. Given that detecting GPT2-generated deepfakes was identified as the most challenging aspect of the task, this accomplishment is considered a success in our research.

6 Conclusions

We have stated research questions which defined the experiments. We ~~managed to investigate them and generally, the answers are positive~~ successfully investigated them and drew conclusions. We

Table 1: Table with 10 best approaches with respect to balanced accuracy

model	dataset	ba	f1	precision	recall	model-dataset
0.8835 <i>ROBERTA (TweepFake)</i>	0.8821 <i>raw</i>	0.8934 0.896	0.8711 0.897	0.891	0.902	
XLM2	raw	0.8835	0.8821	0.8934	0.8711	
SVC	bert	0.8757	0.8763	0.8729	0.8797	
XLM1	<i>SVC</i> raw	<i>bert</i> 0.8713	0.8786	0.8328	0.9297	XLM1-raw-
DISTIL_BERT0	raw	0.8698	0.8686	0.8773	0.8602	DISTIL_BERT0-raw-
GPT2	raw	0.8671	0.8686	0.8593	0.8781	
LGBM	<i>GPT2</i> bert	<i>raw</i> 0.8561	0.8590	0.8429	0.8758	
DISTIL_BERT1	<i>LGBM</i> raw	<i>bert</i> 0.8554	0.8529	0.8681	0.8383	
XGB	<i>DISTIL_BERT1</i> bert	<i>raw</i> 0.8518	0.8546	0.8395	0.8703	
CharCNN+GRU	<i>XGB</i> lemmatized	<i>bert</i> 0.8408	0.8518	0.7975	0.9141	
LR	<i>CharCNN+GRU</i> bert	<i>lemmatized</i> 0.8393	0.8416	0.8304	0.8531	LR-bert-

Table 2: Table with 10 best approaches with respect to global accuracy

model name	TWEET CREATOR (CATEGORY)				
	ALL	GPT2	HUMAN	OTHERS	RNN
<i>ROBERTA_FT (TweepFake)</i>	0.896	0.74	0.89	0.95	1.00
XLM2_raw	0.8835	0.6953	0.8959	0.9153	0.9830
SVC_bert_embeddings	0.8757	0.6927	0.8717	0.9442	0.9782
XLM1_raw	0.8714	0.8307	0.8130	0.9607	0.9854
DistilBERT0_raw	0.8698	0.6589	0.8795	0.9112	0.9879
GPT2_raw	0.8671	0.6693	0.8560	0.9587	0.9782
LGBM_bert_embeddings	0.8561	0.6745	0.8365	0.9483	0.9782
DistilBERT1_raw	0.8554	0.6849	0.8725	0.8471	0.9709
XGB_bert_embeddings	0.8518	0.6562	0.8333	0.9525	0.9733
CharCNN_GRU_lemmatized	0.8409	0.7760	0.7676	0.9628	0.9854
LR_bert_embeddings	0.8393	0.6380	0.8255	0.9236	0.9709

discovered that the use of emoticons, mentions, and URLs is a significant indicator, and in most cases, it suggests that the tweet was generated by a human. We improved the results obtained on GPT-2 generated text by a considerable margin (from 74% to 83% accuracy) without a drop in global accuracy. ~~We can say we (TweepFake: 89%, Our model: 87%).~~ We successfully built a well-performing deepfake detection algorithm. We also found out that effective embeddings are created by high-level deep learning models like BERT. Finally, we recognized the basic patterns of artificially generated texts, i.e. lack of emoticons, mentions, misspells and URLs.

Future works

~~The~~ A good direction of the research is to include more data. Especially datasets generated by more advanced algorithms or relevant to broader corpora than just tweets. An example of such a dataset, which we considered using, is the HC3-English datasets (Guo et al., 2023). It is a group of datasets from different sources, in which for each question, at least one human and ChatGPT3.5 answer is provided.

The questions and answers by the human experts come mainly from publicly available question-answering datasets. There is also an additional source in which Wikipedia is treated as a human expert who is asked questions based on concepts in crawled data.

Another interesting thing would be to extend the current TweepFake dataset with ~~more examples and use more deepfakes generated by more~~ state-of-the-art models ~~to generate fakes~~ like GPT-3.5 (Brown and et al., 2020), GPT-4 (OpenAI et al., 2023), or Gemini (lastly announced Google solution). Unfortunately, with the change in X API policy (formerly known as Twitter), it is much harder because there is no free plan to read tweets, and paid plans are much more restricted than earlier (3,000 reads per month on basic plan). Before that, you could read 300,000 tweets per month on a free plan and even 10,000,000 after applying for a scientific plan. That said, the existing examples with proper prompt engineering can be used with state-of-the-art models to generate even more challenging deepfakes.

Another thing which could be checked with TweepFake is trying to fine-tune more state-of-

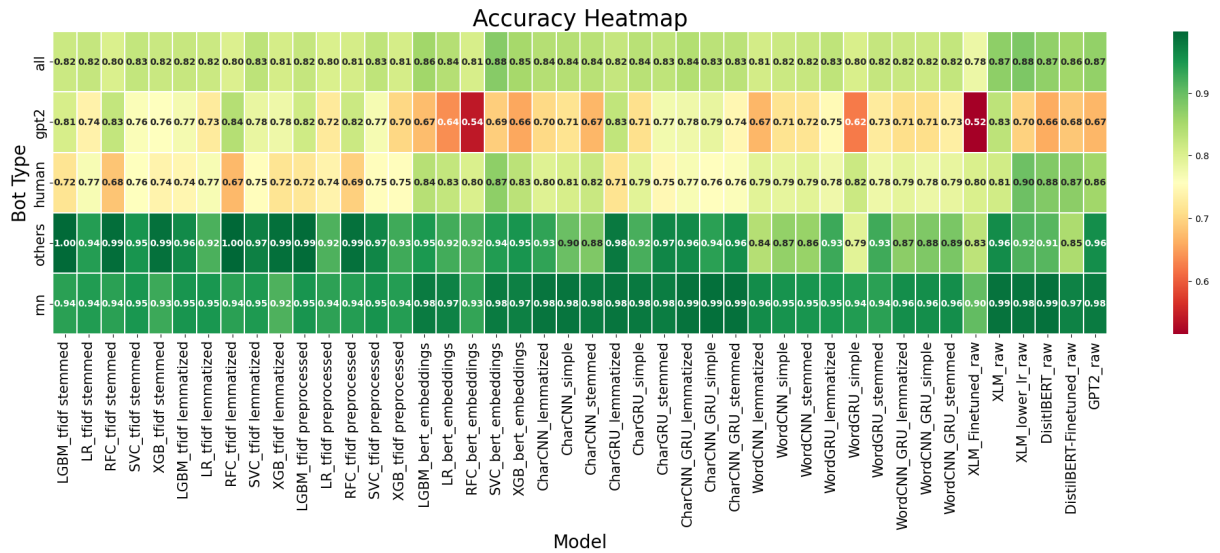


Figure 2: Accuracy for class type and every model

the-art and bigger LLM like LLaMA (Touvron et al., 2023) or any of the earlier mentioned solutions.

Table 3: Contribution table (code in this table adheres to the number within the name of notebook)

person	task	t[h]
Adam Frej	Introduction	3
	Related works (research)	9
	Concept and work plan	4
	Rebuttal I	1
	review I Team 4	2
	Approach	5
	Conclusions	1
	code: eda preprocessing	4
Adrian Kamiński	Related works	6
	Datasets (research)	1.5
	Concept and work plan	4
	Rebuttal I	1
	review I Team 4	2
	Approach	3
	Results	3
	Conclusions	1
	code 01: GPT-2 download, eda	2
	code 41-43,49: transformers	10
	code 91: results	3
Piotr Marciniak	Introduction	4
	Datasets (writing & research)	3
	Concept and work plan	3
	Approach	3

person	task	t[h]
	Rebuttal I	2
	review I Team 10	2.5
	code 01-12	10
	code 21-23 (modifications)	1.5
	code 31-33	5
	code 90	3
Szymon	Related works (research)	4
Szmajdziński	Datasets (research)	3
	Concept	3
	Abstract	1
	review I Team 10	2.5
	code 21-23: Deep Learning Models	10
	Results	5
	Conclusions	1

References

- [Adelani et al.2019] David Ifeoluwa Adelani, Haotian Mai, Fuming Fang, Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. 2019. Generating sentiment-preserving fake online reviews using neural language models and their human- and machine-based detection. *CoRR*, abs/1907.09177.
- [Akiba et al.2019] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [Bakhtin et al.2019] Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng, Marc’Aurelio Ranzato, and Arthur Szlam. 2019. Real or fake? learning to discriminate machine from human generated text. *CoRR*, abs/1906.03351.
- [Barrus2019] Tyler Barrus. 2019. pypellchecker. <https://pypi.org/project/pypellchecker/>. accessed: 09.12.2023. [Online].
- [Bird et al.2009] Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python*. O’Reilly Media, Inc., 1st edition.
- [Brown and et al.2020] Tom B. Brown and Benjamin Mann et al. 2020. Language models are few-shot learners.
- [Chen and Guestrin2016] Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, pages 785–794, New York, NY, USA. ACM.
- [Conneau et al.2019] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116.
- [Cortes and Vapnik1995] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.
- [Cox1958] David R Cox. 1958. The regression analysis of binary sequences. *Journal of the Royal Statistical Society: Series B (Methodological)*, 20(2):215–232.
- [Devlin et al.2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.
- [Europol Innovation Lab2024] Europol Innovation Lab. 2024. Facing reality? Law enforcement and the challenge of deepfakes. *Europol*. Accessed on 28 January 2024.
- [Fagni et al.2021] Tiziano Fagni, Fabrizio Falchi, Margherita Gambini, Antonio Martella, and Maurizio Tesconi. 2021. Tweepfake: About detecting deepfake tweets. *PLOS ONE*, 16(5):1–16, 05.
- [Fernandez2024] Ray Fernandez. 2024. YouTube Deepfake and AI Crypto Scams Take \$600K With “Double Your Money!” Promise. *Techopedia*. Accessed on 28 January 2024.
- [Gehrmann et al.2019] Sebastian Gehrmann, Hendrik Strobelt, and Alexander M. Rush. 2019. Gltr: Statistical detection and visualization of generated text.
- [Guo et al.2023] Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection.
- [He and McAuley2016] Ruining He and Julian J. McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. *CoRR*, abs/1602.01585.
- [Ho1995] Tin Kam Ho. 1995. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE.
- [Honnibal et al.2020] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python.
- [Ippolito et al.2020] Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. Automatic detection of generated text is easiest when humans are fooled. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors,

- Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822, Online, July. Association for Computational Linguistics.
- [Ke et al.2017] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30:3146–3154.
- [Liu et al.2019] 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.
- [Lynn2024] Bryan Lynn. 2024. ‘Deepfake’ of Biden’s Voice Called Early Example of US Election Disinformation. *Voice of America*. Accessed on 28 January 2024.
- [Milmo2023] Dan Milmo. 2023. Doctored Sunak picture is just latest in string of political deepfakes. *The Guardian*. Accessed on 28 January 2024.
- [OpenAI et al.2023] OpenAI, :, Josh Achiam, Steven Adler, and et al. 2023. Gpt-4 technical report.
- [Parsons2022] Jeff Parsons. 2022. Ukraine warns Russia may deploy deepfakes of Volodymyr Zelensky surrendering. *Metro*. Accessed on 28 January 2024.
- [Radford et al.2019] Alec Radford, Jeff Wu, Rewon Child, D. Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- [Rahman-Jones2024] Imran Rahman-Jones. 2024. Taylor Swift deepfakes spark calls in Congress for new legislation. *BBC*. Accessed on 28 January 2024.
- [Sanh et al.2020] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.
- [Sebastiani2002] Fabrizio Sebastiani. 2002. Machine learning in automated text categorization. *ACM Comput. Surv.*, 34(1):1–47, mar.
- [Sennrich et al.2016] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units.
- [Solaiman et al.2019] Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, Gretchen Krueger, Jong Wook Kim, Sarah Kreps, et al. 2019. Release strategies and the social impacts of language models. *arXiv preprint arXiv:1908.09203*.
- [Touvron et al.2023] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.
- [Vincent2018] James Vincent. 2018. Why we need a better definition of ‘deepfake’ / let’s not make deepfakes the next fake news. <https://www.theverge.com/2018/5/22/17380306/deepfake-definition-ai-manipulation-fake-news>, May.
- [Yang et al.2020] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2020. Xlnet: Generalized autoregressive pretraining for language understanding.
- [Zellers et al.2019a] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019a. Defending against neural fake news. *CoRR*, abs/1905.12616.
- [Zellers et al.2019b] Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019b. Defending against neural fake news. *CoRR*, abs/1905.12616.
- [Zhang et al.2015] Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *CoRR*, abs/1509.01626.

Appendix

Runtime parameters

Average runtime for each approach (in the following format hour : minute : second).

- 00:05:12 CharCNN (33 epochs)
- 00:28:33 CharGRU (24 epochs)
- 00:50:04 CharCNN+GRU (45 epochs)
- 00:01:43 WordCNN (12 epochs)
- 00:08:38 WordGRU (14 epochs)
- 00:06:44 WordCNN+GRU (12 epochs)
- 1:15:41 XML0 (10 epochs)
- 1:31:45 XML1 (4 epochs)
- 1:54:42 XML2 (5 epochs)
- 53:44 DISTIL_BERT0 (5 epochs)
- 1:19:33 (GPT2) + 53:46 (TweepFake) DISTIL_BERT1 (10, 500 GPT2 samples + 5 epochs)
- 2:49:42 GPT2 (7 epochs)

Infrastructure

All experiments were concluded on a single-machine setup.

We used a machine with an Intel Core i5-13400F processor and Nvidia RTX3060 GPU for all experiments.

Reproducibility checklist

Overall results:

- MODEL DESCRIPTION – A clear description of the mathematical setting, algorithm, and/or model

Stated in Section 4.

- LINK TO CODE – A link to a downloadable source code, with specification of all dependencies, including external libraries.

Link to source code with all information about reproducing the results – github.com/grant-TraDA/NLP-2023W/tree/main/15. Deepfake tweets detection.

- INFRASTRUCTURE – A description of the computing infrastructure used.

Stated in 3

- RUNTIME PARAMETERS – Average runtime for each approach

Stated in Section 3 (Appendix).

- PARAMETERS – The number of parameters in each model

Not reported.

- VALIDATION PERFORMANCE – Corresponding validation performance for each reported test result
Stated in Section 5.
- METRICS – Explanation of evaluation metrics used, with links to code
Stated in Section 4.

Multiple Experiments:

- NO TRAINING EVAL RUNS – The exact number of training and evaluation runs
Stated in Section 4 and Section 3 (Appendix).
- HYPER BOUND – Bounds for each hyperparameter
Stated in Section 4 and Table 4.
- HYPER BEST CONFIG – Hyperparameter configurations for best-performing models
Stated in Section 4.
- HYPER SEARCH – Number of hyperparameter search trials
Stated in Section 4.
- HYPER METHOD – The method of choosing hyperparameter values (e.g., uniform sampling, manual tuning, etc.) and the criterion used to select among them (e.g., accuracy)
Stated in Section 4.
- EXPECTED PERF – Summary statistics of the results (e.g., mean, variance, error bars, etc.)
Stated in Section 5.

Datasets – utilized in the experiments and/or the created ones:

- DATA STATS – Relevant statistics, such as the number of examples
Stated in Section 3 and Section 4 as well as in EDA (Images in Appendix/notebooks in repository with source code).
- DATA SPLIT – Details of train/validation/test splits
Stated in Section 4.
- DATA PROCESSING – Explanation of any data that were excluded and all preprocessing steps
Stated in Section 4.2.
- DATA DOWNLOAD – A link to a downloadable version of the data
Stated explicitly in the repository README file and implicitly as a reference to a paper that created the dataset.
- NEW DATA DESCRIPTION – For new data collected, a complete description of the data collection process, such as instructions to annotators and methods for quality control.
Not applicable.
- DATA LANGUAGES – For natural language data, the name of the language(s).
English.

Rebuttal from the second part of the project

Team 5

- The hypotheses mentioned in the report, despite being interesting, were not significantly impactful for the research.

We agree that the hypotheses did not prove to be impactful for our research. This is one of the points discussed in our work. Generally, it is a challenging task to state good hypotheses at the beginning of the experiments. Moreover, we believe that despite not giving directly good results, they are still useful for our work and help draw further research direction.

- The authors could explore additional datasets for the test.

We agree that additional datasets would be helpful for further research, which is discussed in the Future works. That being said, our experiment presented in 4.5 did not improve results but worsened. We hope this work provides sufficient grounds for the research.

Team 14

- Relatively low level of English. There are multiple grammar mistakes, which could have been omitted with the usage of tools, such as Grammarly. Some sentences don't make sense and seem to end prematurely.

We believe we improved the quality of the article's language. We redacted the text and ensured its consistency.

- The conclusions seem to be very general, and in some cases don't point out particular gains.

Particular gains and specific pros and cons of our research are described in the Results chapter. The conclusions are purposefully general, summing up the work.

- The description of the research methodology is extremely chaotic, and it is hard to understand and visualize the experimental pipeline.

Poor, and chaotic structure of the paper, which makes it harder to fully focus on particular parts of the study such as preprocessing, model specification, or the analysis of the outcomes.

We agree that the original version of the paper was chaotic and not intuitive to read. We added more sections and subsection indicators, as well as restructured particular paragraphs. We believe the work should be easier to comprehend, with a special focus on research methodology.

- Why did you discard 2 datasets not being TweepFake? Please put this information in the paper.

We rephrased the datasets' descriptions to make our approach more straightforward. We discarded only one dataset, which is put in section 6 as consideration for usage in future works.

- Suggestions and Missing references

We are grateful for the provided tips and ideas. We considered them and incorporated some of them into our work. The pointed, small mistakes helped to make the paper cleaner.

Table 4: The table presenting model and encoder hyperparameters which were optimized by the optuna package (Integer - means that values are integers from the given range, Float - floats from the given range, Categorical - category from the given list, log=True - means that the values are drawn from log-uniform distribution between given range)

model / encoder	hyperparameter	value space
LGBM	Boosting Type	Categorical([gbdt, dart])
	max_depth	Integer(1, 15)
	n_estimators	Integer(10, 500, log=True)
	subsample	Float(0.6, 1)
XGB	booster	Categorical([gbtree, dart])
	max_depth	Integer(1, 15)
	n_estimators	Integer(10, 500, log=True)
	subsample	Float(0.6, 1)
RF	max_depth	Integer(1, 15)
	n_estimators	Integer(10, 500, log=True)
	criterion	Categorical([gini, entropy, log_loss])
	min_samples_split	Float(0.01, 0.1)
SVC	kernel	Categorical([linear, poly, rbf, sigmoid])
	C	Float(1e-2, 1e2, log=True)
Logistic regression	penalty	Categorical([l1, l2])
	C	Float(1e-2, 1e2, log=True)
TF-IDF	ngram_range	Categorical([(1, 1), (1, 2), (1, 3)])
	max_features	Integer(1000, 10000, log=True)
	max_df	Float(0.8, 1)
	min_df	Float(0.0, 0.2)

Table 5: Table with results of all experiments with ~~with~~ more specific metrics

balanced_accuracy	f1_score	precision	recall	model	dataset
0.8835	0.8821	0.8934	0.8711	XLM2	raw
0.8757	0.8763	0.8729	0.8797	SVC	bert_embeddings
0.8713	0.8786	0.8328	0.9297	XLM1	raw
0.8698	0.8686	0.8773	0.8602	DISTIL_BERT0	raw
0.8671	0.8686	0.8593	0.8781	GPT2	raw
0.8561	0.8590	0.8429	0.8758	LGBM	bert_embeddings
0.8554	0.8529	0.8681	0.8383	DISTIL_BERT1	raw
0.8518	0.8546	0.8395	0.8703	XGB	bert_embeddings
0.8408	0.8518	0.7975	0.9141	CharCNN+GRU	lemmatized
0.8393	0.8416	0.8304	0.8531	LR	bert_embeddings
0.8385	0.8451	0.8125	0.8805	CharCNN	lemmatized
0.8381	0.8432	0.8184	0.8695	CharCNN	simple
0.8354	0.8420	0.8101	0.8766	CharGRU	simple
0.8350	0.8374	0.8260	0.8492	CharCNN	stemmed
0.8330	0.8442	0.7919	0.9039	CharCNN+GRU	stemmed

Continued on next page

balanced_accuracy	f1_score	precision	recall	model	dataset
0.8322	0.8444	0.7881	0.9094	CharCNN+GRU	simple
0.8311	0.8433	0.7873	0.9078	SVC	lemmatized
0.8291	0.8423	0.7827	0.9117	CharGRU	stemmed
0.8287	0.8373	0.7982	0.8805	WordGRU	lemmatized
0.8287	0.8407	0.7864	0.9031	SVC	simple
0.8260	0.8373	0.7869	0.8945	SVC	stemmed
0.8248	0.8315	0.8019	0.8633	WordCNN+GRU	stemmed
0.8244	0.8323	0.7974	0.8703	WordGRU	stemmed
0.8244	0.8340	0.7916	0.8812	LR	stemmed
0.8244	0.8421	0.7658	0.9352	CharGRU	lemmatized
0.8236	0.8402	0.7686	0.9266	LGBM	simple
0.8225	0.8356	0.7787	0.9016	LGBM	lemmatized
0.8209	0.8257	0.8049	0.8477	WordCNN	stemmed
0.8201	0.8256	0.8019	0.8508	WordCNN+GRU	lemmatized
0.8189	0.8325	0.7751	0.8992	XGB	stemmed
0.8186	0.8241	0.8004	0.8492	WordCNN	simple
0.8186	0.8276	0.7890	0.8703	LR	lemmatized
0.8185	0.8356	0.7646	0.9211	LGBM	stemmed
0.8178	0.8244	0.7962	0.8547	WordCNN+GRU	simple
0.8100	0.8134	0.7998	0.8273	WordCNN	lemmatized
0.8092	0.8256	0.7609	0.9023	XGB	lemmatized
0.8084	0.8185	0.7782	0.8633	XGB	simple
0.8080	0.8278	0.7511	0.9219	RF	simple
0.8061	0.8069	0.8043	0.8094	RF	bert_embeddings
0.8045	0.8019	0.8135	0.7906	WordGRU	simple
0.8037	0.8251	0.7447	0.9250	RF	stemmed
0.8021	0.8140	0.7688	0.8648	LR	simple
0.8013	0.8243	0.7395	0.9313	RF	lemmatized
0.7811	0.7764	0.7941	0.7594	XLM0	raw

Table 6: Table with results of all experiments with split distinguishing the bot type responsible for the creation of twitter posts.

bot_type	all	gpt2	human	others	rnn
model_name					
XLM2_raw	0.8835	0.6953	0.8959	0.9153	0.9830
SVC_bert_embeddings	0.8757	0.6927	0.8717	0.9442	0.9782
XLM1_raw	0.8714	0.8307	0.8130	0.9607	0.9854
DisitlBERT0_raw	0.8698	0.6589	0.8795	0.9112	0.9879
GPT2_raw	0.8671	0.6693	0.8560	0.9587	0.9782
LGBM_bert_embeddings	0.8561	0.6745	0.8365	0.9483	0.9782
DistilBERT1_raw	0.8554	0.6849	0.8725	0.8471	0.9709
XGB_bert_embeddings	0.8518	0.6562	0.8333	0.9525	0.9733
CharCNN_GRU_lemmatized	0.8409	0.7760	0.7676	0.9628	0.9854
LR_bert_embeddings	0.8393	0.6380	0.8255	0.9236	0.9709
CharCNN_lemmatized	0.8385	0.7031	0.7966	0.9339	0.9830

Continued on next page

bot_type model_name	all	gpt2	human	others	rnn
CharCNN_simple	0.8382	0.7135	0.8067	0.8967	0.9830
CharGRU_simple	0.8354	0.7109	0.7942	0.9194	0.9806
CharCNN_stemmed	0.8350	0.6693	0.8208	0.8822	0.9782
CharCNN_GRU_stemmed	0.8331	0.7448	0.7621	0.9587	0.9879
CharCNN_GRU_simple	0.8323	0.7865	0.7551	0.9380	0.9903
SVC_tfidf lemmatized	0.8311	0.7839	0.7543	0.9731	0.9466
CharGRU_stemmed	0.8292	0.7734	0.7465	0.9669	0.9757
WordGRU_lemmatized	0.8288	0.7474	0.7770	0.9277	0.9490
SVC_tfidf preprocessed	0.8288	0.7682	0.7543	0.9669	0.9539
SVC_tfidf stemmed	0.8260	0.7630	0.7574	0.9525	0.9490
WordCNN_GRU_stemmed	0.8249	0.7318	0.7864	0.8864	0.9587
WordGRU_stemmed	0.8245	0.7266	0.7786	0.9256	0.9393
LR_tfidf stemmed	0.8245	0.7370	0.7676	0.9421	0.9442
CharGRU_lemmatized	0.8245	0.8281	0.7136	0.9793	0.9830
LGBM_tfidf preprocessed	0.8237	0.8151	0.7207	0.9938	0.9515
LGBM_tfidf lemmatized	0.8225	0.7708	0.7433	0.9628	0.9515
WordCNN_stemmed	0.8210	0.7161	0.7942	0.8636	0.9515
WordCNN_GRU_lemmatized	0.8202	0.7083	0.7895	0.8719	0.9587
XGB_tfidf stemmed	0.8190	0.7578	0.7387	0.9876	0.9272
LGBM_tfidf stemmed	0.8186	0.8073	0.7160	0.9959	0.9393
LR_tfidf lemmatized	0.8186	0.7266	0.7668	0.9174	0.9490
WordCNN_simple	0.8186	0.7135	0.7879	0.8678	0.9539
WordCNN_GRU_simple	0.8178	0.7083	0.7809	0.8822	0.9587
WordCNN_lemmatized	0.8100	0.6693	0.7926	0.8388	0.9612
XGB_tfidf lemmatized	0.8092	0.7786	0.7160	0.9897	0.9150
XGB_tfidf preprocessed	0.8084	0.7005	0.7535	0.9256	0.9417
RFC_tfidf preprocessed	0.8081	0.8177	0.6941	0.9876	0.9417
RFC_bert_embeddings	0.8061	0.5417	0.8028	0.9174	0.9320
WordGRU_simple	0.8045	0.6224	0.8185	0.7934	0.9442
RFC_tfidf stemmed	0.8038	0.8255	0.6823	0.9917	0.9393
LR_tfidf preprocessed	0.8022	0.7240	0.7394	0.9153	0.9369
RFC_tfidf lemmatized	0.8014	0.8385	0.6714	0.9979	0.9393
XML0_raw	0.7811	0.5156	0.8028	0.8306	0.9029

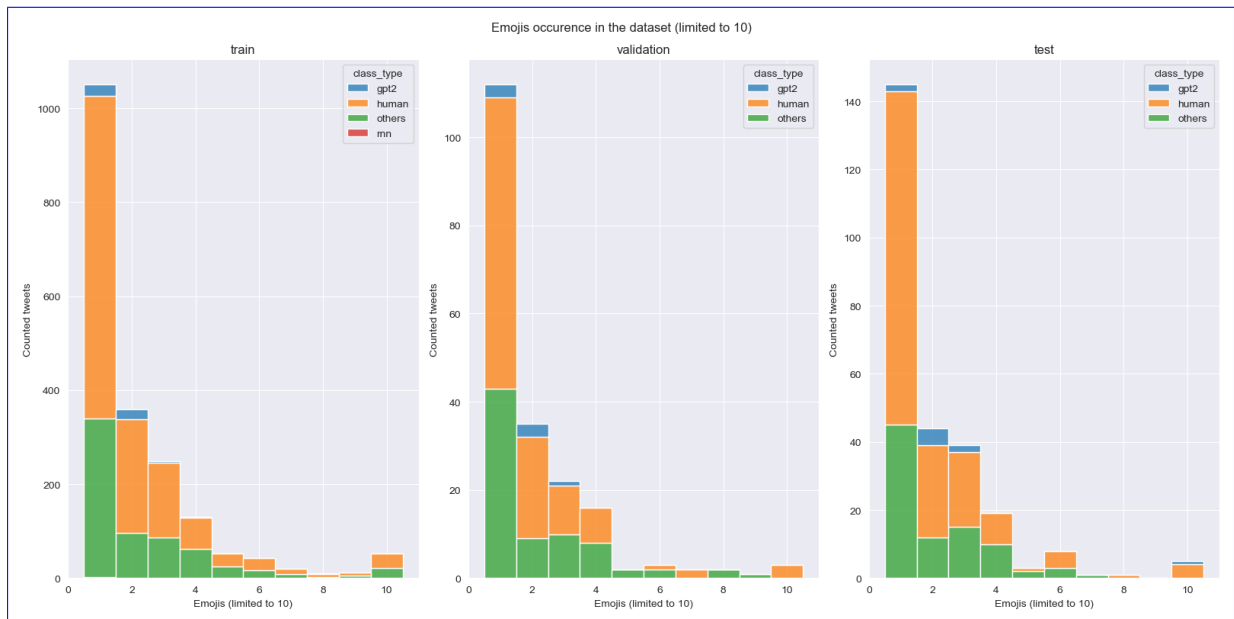


Figure 3: [Emoticons analysis](#)

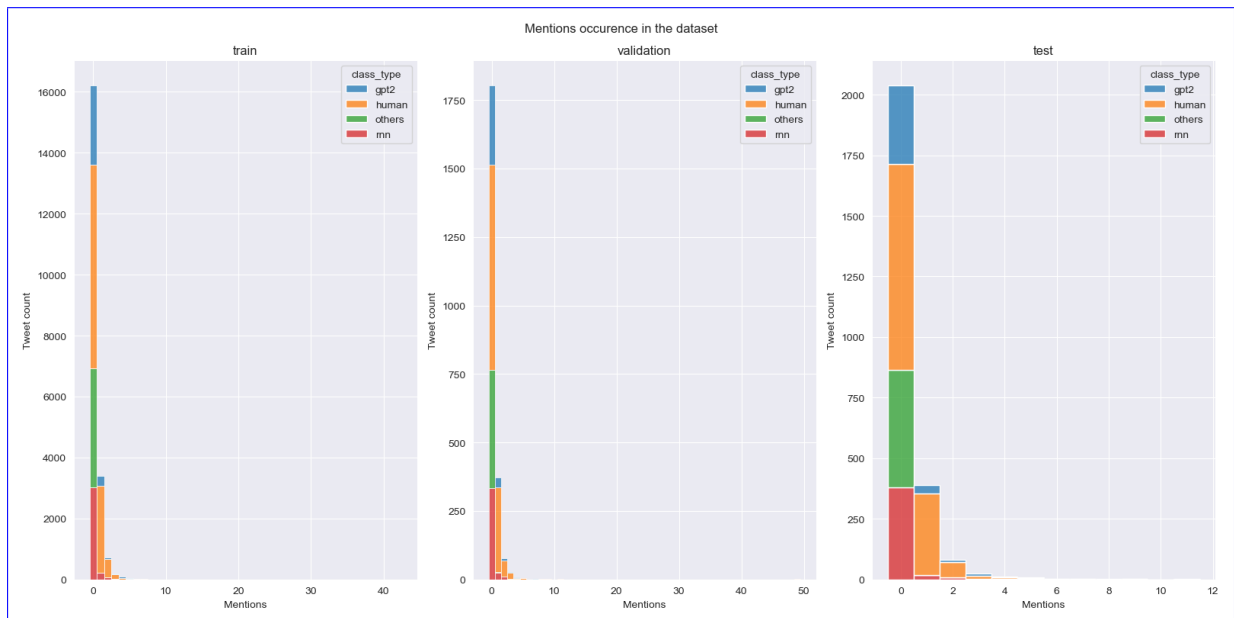


Figure 4: [Mentions analysis](#)

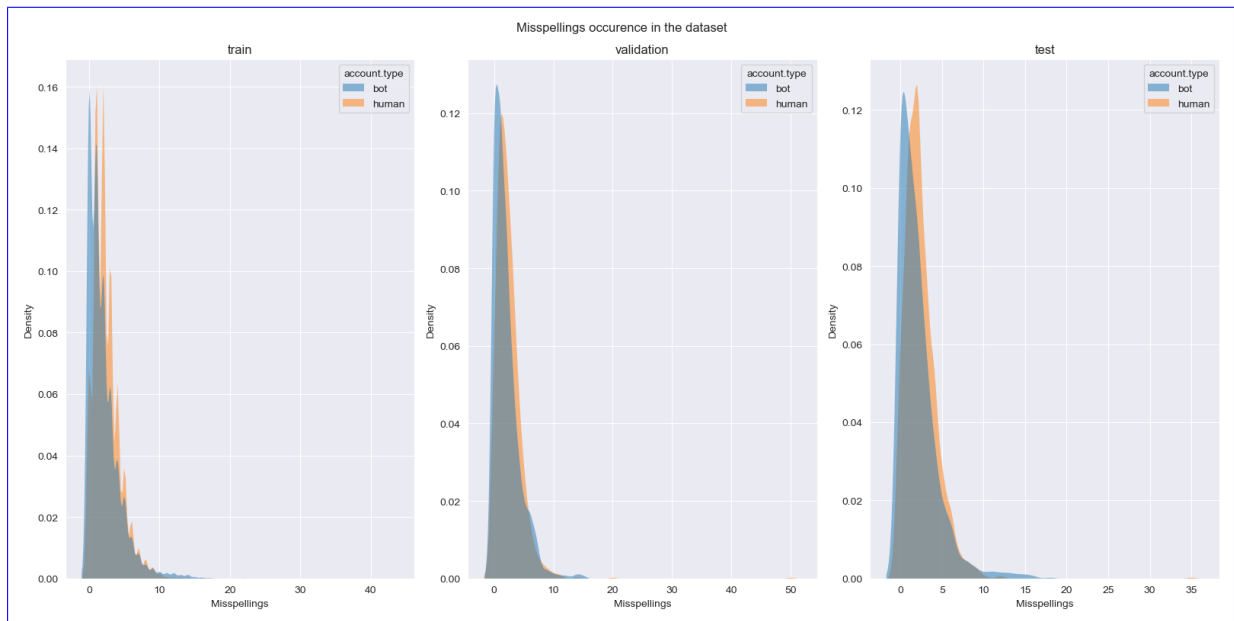


Figure 5: Misspellings analysis

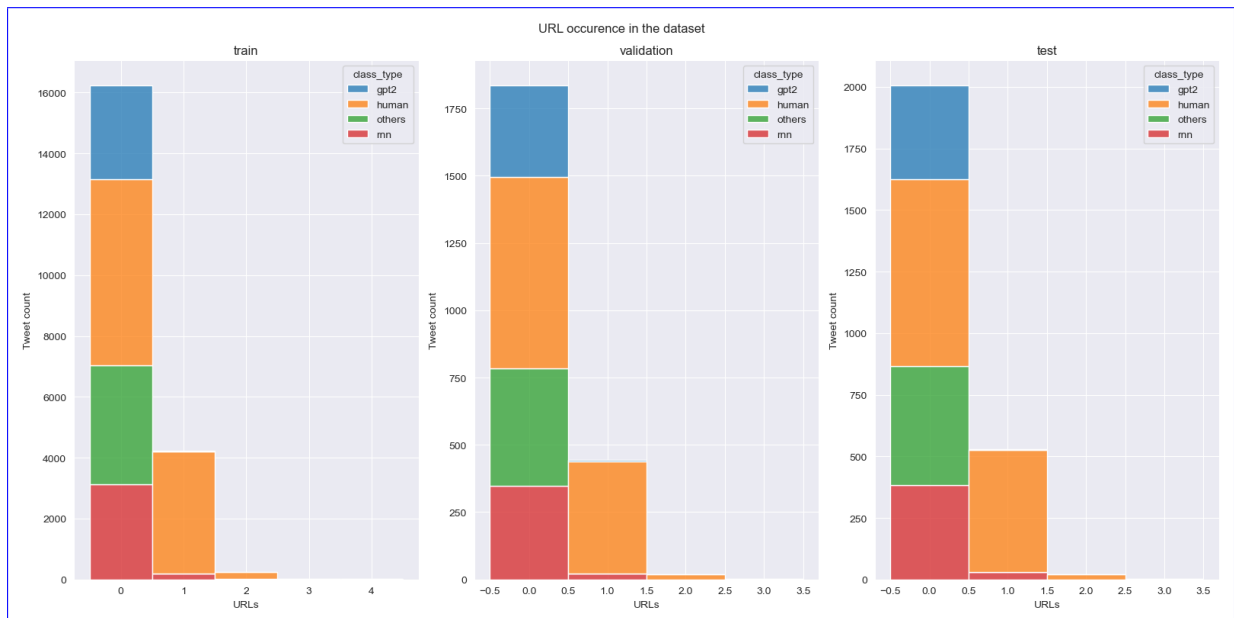


Figure 6: URLs analysis