




Article

Are Strong Baselines Enough? False News Detection with Machine Learning

Lara Aslan ^{1,†}, Michal Ptaszynski ^{1,*,†}  and Jukka Jauhiainen ^{2,*}

¹ Text Information Processing Laboratory, Kitami Institute of Technology, Kitami 090-8507, Japan; m3245380016@std.kitami-it.ac.jp

² School of Information Technology, Oulu University of Applied Sciences, 90570 Oulu, Finland

* Correspondence: michal@mail.kitami-it.ac.jp (M.P.); jukka.jauhiainen@oamk.fi (J.J.)

† These authors contributed equally to this work.

Abstract: False news refers to false, fake, or misleading information presented as real news. In recent years, there has been a noticeable increase in false news on the Internet. The goal of this paper was to study the automatic detection of such false news using machine learning and natural language processing techniques and to determine which techniques work the most effectively. This article first studies what constitutes false news and how it differs from other types of misleading information. We also study the results achieved by other researchers on the same topic. After building a foundation to understand false news and the various ways of automatically detecting it, this article provides its own experiments. These experiments were carried out on four different datasets, one that was made just for this article, using 10 different machine learning methods. The results of this article were satisfactory and provided answers to the original research questions set up at the beginning of this article. This article could determine from the experiments that passive aggressive algorithms, support vector machines, and random forests are the most efficient methods for automatic false news detection. This article also concluded that more complex experiments, such as using multiple levels of identifying false news or detecting computer-generated false news, require more complex machine learning models.

Keywords: machine learning; natural language processing; false news detection; artificial intelligence; ChatGPT



Citation: Aslan, L.; Ptaszynski, M.; Jauhiainen, J. Are Strong Baselines Enough? False News Detection with Machine Learning. *Future Internet* **2024**, *16*, 322. <https://doi.org/10.3390/fi16090322>

Academic Editor: Filipe Portela

Received: 13 August 2024

Revised: 26 August 2024

Accepted: 27 August 2024

Published: 5 September 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

False news can be described as false or misleading information created to be widely shared for different purposes, such as to generate revenue, promote, or affect the opinions of a target group before a big event. Recent years have seen the spread of false news, especially on the Internet [1–3]. This issue needs to be addressed, as the growing spread of false news threatens journalism, can cause political turmoil, and can negatively impact people’s daily lives [4].

While the internet already has many websites that professionally perform fact-checking on news articles, such as PolitiFact (<https://www.politifact.com/>, accessed on 12 August 2024), FactCheck.org (<https://www.factcheck.org/>, accessed on 12 August 2024), FactChecker (<https://www.washingtonpost.com/news/fact-checker/>, accessed on 12 August 2024), Snopes (<https://www.snopes.com/>, accessed on 12 August 2024), The Reporters’ Lab (<https://reporterslab.org/>, accessed on 12 August 2024), or FaktaBaari (<https://www.faktaaari.fi/fakta/>, accessed on 12 August 2024), in today’s world, websites like these managed by human volunteers are not enough. False news or false information is becoming harder to detect as the ways to create it have become more advanced and accessible [5]. For example, ChatGPT (<https://chat.openai.com/>, accessed on 12 August 2024) can just in a few seconds generate an article, a review, or a statement that can be false but still look believable [6–8]. This phenomena can be also observed in Figure 1. With the

help of social media, these pieces of false information can spread to millions of people in just a few seconds. For this reason, it is important to include at least some automation in the fact-checking process.

ChatGPT ▾

Lets play creative word games. Write a piece of text that would look like an article about Vladimir Putin being friends with Bigfoot.



In a shocking revelation that is sure to shake up the world of international diplomacy, it has been reported that Russian President Vladimir Putin has struck up an unlikely friendship with none other than Bigfoot.

According to sources close to the Kremlin, the two unlikely companions first met in the remote forests of Siberia, where Putin was reportedly on a hunting trip. It is said that the Russian leader was initially startled by the appearance of the legendary

Figure 1. An example of succeeding in tricking ChatGPT into writing a false news article.

In this article, we study automated false news detection. We review the previous work related to this subject and analyze such data and classification experiments using machine learning algorithms. Specifically, we test several classification models using machine learning and natural language processing techniques to find the best way to automatically detect when a news article or a piece of information is true or false and analyze the language commonly used in news articles that could be considered false. We also experimented on how the developed classifiers perform when given articles written by a human compared to those generated automatically.

The usage of generative large language models, including the recently popular ChatGPT makes it easier to write false content in large quantities. Thus, the inclusion of data generated with ChatGPT is important when it comes to false news detection.

Most fact-checking services or previous research conducted in automated false news detection focuses mainly on political issues. False news, however, can be spread about any topic, including health-related ones, such as COVID-19 [9]. Therefore, it is important to include other topics in the scope of analysis as well.

Much of the previous research conducted in automated false news detection often only uses one dataset [10–12]. The dependence on particular datasets might cause issues if these solutions, which have only been evaluated on a single dataset, are used to detect more diverse or real-world data. To reduce the potential issues that using a single dataset could cause, our research includes four distinct datasets.

We also have to take into consideration that the scope of deceptive information in today's age is large, and often, different definitions of different types of deceptive information are used indiscriminately with each other. This creates confusion and uncertainty when trying to differentiate one piece of deceptive information from another.

Different types of deceptive information include rumors, hoaxes, false news, false reviews, satires, urban legends, and propaganda, among many more. It is important to take into consideration that sometimes these different types of deceptive information can intertwine with each other. For example, a false news article can be satirical. To understand better the differences between these different types of deceptive information we summarized all known types of deceptive information so far with their short definitions in Table 1.

Table 1. Descriptions of different types of deceptive information types.

Type	Description
Rumor	Quickly spreading story or news that can be true or invented.
A hoax	A deceptive piece of information used to trick people into believing in it.
False news	False or misleading information is presented as news. Used to be widely shared for influencing purposes.
False reviews	A review that is not an actual consumer’s opinion or does not reflect the actual opinion of a consumer. Often used to manipulate a consumer not to buy a certain product.
Satires	A type of parody where content is presented with irony or humor. Often used to criticize events, people, etc.
Urban legends	A false story that is circulated between people as true. Usually humorous, horrifying, or cautionary.
Propaganda	Information that is usually biased or misleading. It is used to promote a political cause or a point of view.

In this research, as previously pointed out, we mainly focus on detecting false news. However, even false news can be divided into sub-categories where each different sub-category has a specific purpose and a different impact level. This is why various guides and taxonomies have been created to differentiate and understand different types of false news [13–15].

Wardle [14] proposed a taxonomy of misinformation and disinformation that divided false news into seven different types of mis- and disinformation and provided a description of the intended harm that these different types of false news are causing. Our research mainly focuses on misleading content and fabricated content. This taxonomy is shown in Table 2.

Table 2. “7 types of mis- and disinformation”. A taxonomy created by Wardle [14].

Type	Description
Satire or parody	No intention to cause harm but has the potential to fool.
Misleading Content	Misleading use of information to frame an issue or individual.
Imposter Content	When genuine sources are impersonated.
Fabricated Content	New content is 100% fake, designed to deceive and do harm.
False Connection	When headlines, visuals or captions do not support the content, also known as clickbait.
False Context	When genuine content is shared with false contextual information.
Manipulated Content	When genuine information or imagery is manipulated to deceive.

To understand false news, we also need to better understand the definition of news that can be classified as real. For news to be considered legitimate or real, it needs to meet certain journalistic standards. This standard usually means that the news is neutral, uses the right sources, and is factual based on the information available at the time. Chong [16], in their research on misinformation, considered news to be legitimate or real if it followed the following characteristics:

1. Presented in a neutral, balanced, and non-inciting manner.
2. Verifiable by an independent source or party within reasonable limits.
3. Accurate and factual, based on the information available or as provided by the source.
4. Comprehensive—with no malicious censorship, modification, or manipulation.

For this research, we look out for these characteristics proposed by Chong [16] to differentiate between real news and false news. For the terminology used in this paper, please refer to Table 3.

Table 3. Glossary of different terms used in the research.

Type	Description
False news	False or misleading information presented as news, often called “fake news” as well [17]. This paper uses the term false news, as it has a less polarizing connotation.
Machine learning (ML)	The use and development of computer systems that can learn and adapt by using algorithms and statistical data to analyze patterns from the given data [18]
Artificial intelligence (AI)	The ability of computers to perform tasks that are usually more associated with intelligent beings [19]
Large language models (LLMs)	An example of generative AI. They can recognize, translate, predict, or generate texts or other forms of content [20]. A good example of LLMs would be ChatGPT.
Deep learning	A subset of ML, is a neural network with three or more layers. The neural networks attempt to simulate the behavior of the human brain [21]
Natural language processing (NLP)	The application of computational techniques to analyze and synthesize natural language and speech [22]

The two main research questions posed by this research were as follows:

1. What are the best computational methods to use when detecting false news?
2. Will there be a difference in results when using human-generated text and automatically generated text?

The remainder of this article is arranged as follows. In Section 2, we present the previous research performed in the field of false news detection, as well as introduce the datasets used in the research and the applied methods. In Section 3, we introduce the experiments we performed on the datasets we used. In Section 4, we discuss the results of the experiments in more detail, as well as include a minor linguistic analysis of false news. In Section 5, we conclude the research.

2. Materials and Methods

In this section, we present the materials and methods relevant to our research field. This includes the previous research relevant to our experiments, the datasets we used, and the applied methods we implemented for our experiments. Figure 2 presents a flowchart of our research methodology.

2.1. Previous Research

In this sub-section, we go through the previous research studies that are important in automated false news detection and those that are relevant to our research and experiments.

2.1.1. False News Detection

Rubin [23] studied satirical news and how to expose it as false news. The contrast between satirical news and false news is worth noting. Satirical news gives text cues to reveal the false nature of the news, whereas false news tries to convince the reader to believe in it. In their research, Rubin [23] proposed an algorithm based on support vector machines (SVM) with five general features: absurdity, humor, grammar, negative affect, and punctuation to predict satirical news. Their research was very successful, as they achieved a 90% precision and 84% recall in satirical news detection. Rubin’s [23] study effectively uses SVM and important linguistic traits to predict satirical news, achieving

high precision and recall scores. However, the algorithm’s reliance on predetermined cues like absurdity might hinder the algorithm’s adaptability to more complex and evolving forms of satire.

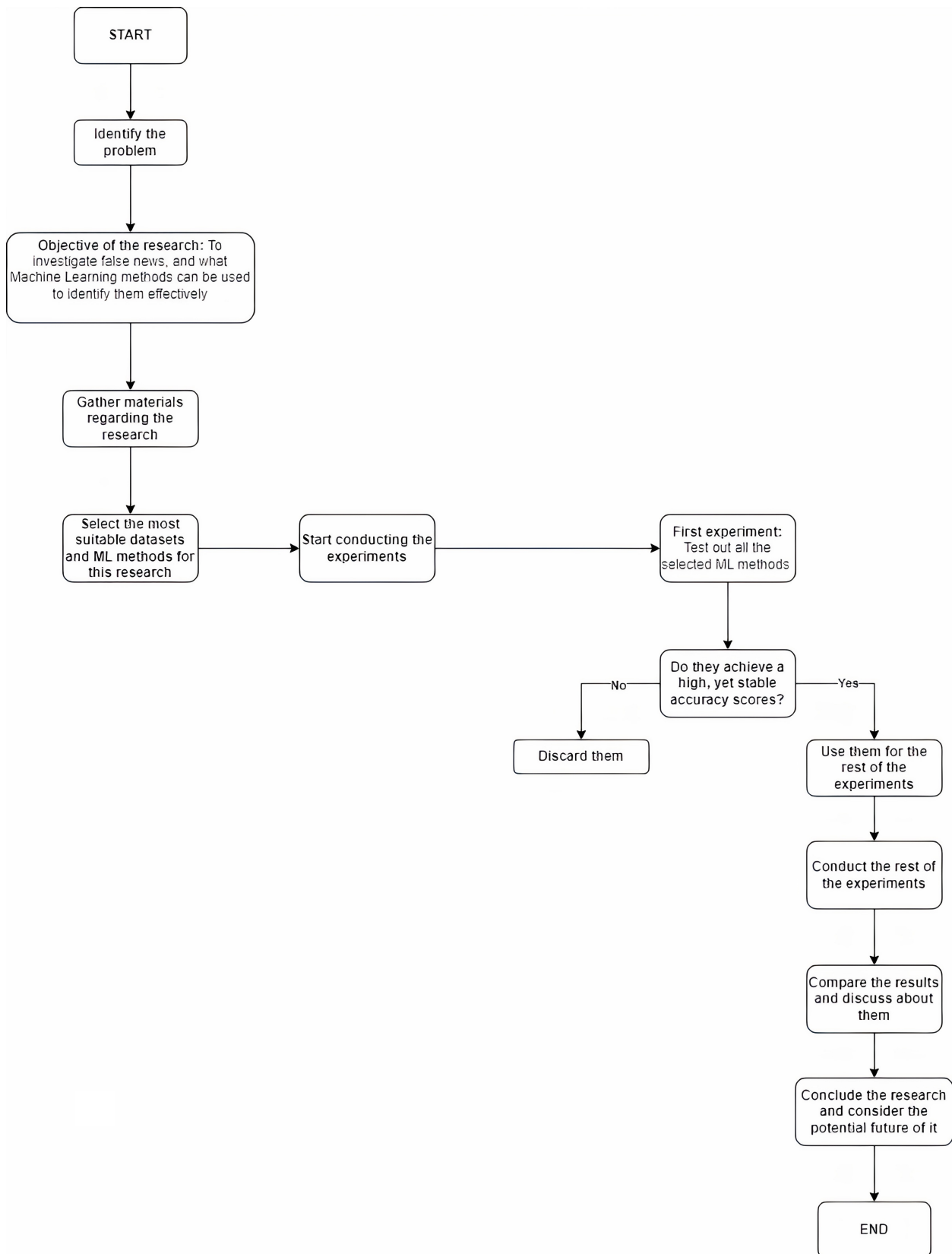


Figure 2. A flowchart of the research methodology.

Thota [11] used deep learning architectures to detect false news. They highlighted that a problem with the majority of false news detectors is that they only use binary classification methods, making them unable to understand the relationship between two pieces of text. In their research, they tackle this problem through stance detection, in which they use deep neural network architecture to predict how similar the headline is to a news article. Their model proved to be successful, as they were able to detect when a news article was false with stance detection with 94.21% accuracy, which outperformed the existing models at the time by 2.5%. While Thota's [11] research addresses the limitations of binary classification in false news detection, their proposed method's focus on headline–article similarity may miss other important aspects of misinformation.

Karimi [24] conducted their research on false news detection with the inclusion of various degrees of “falseness”. They propose a multi-source multi-class fake news detection framework to tackle this problem. This framework combines automated feature extraction, multi-source fusion, and automated degrees of falseness detection into one single model. Their model could differentiate between the different degrees of falseness from the news that they used. They also integrated multiple sources into false news detection, which could help false news detection as multiple sources give a much better context when detecting false news, as opposed to only using the context given by the news article. However, the model's complexity may present challenges for scalability.

Oshikawa [25] studied the potentials and limitations of NLP solutions in false news detection. NLP techniques are perhaps the most common way of analyzing false news, and their study proves that NLP techniques are useful in automatic false news detection. Das [26] additionally studied the task of automatic fact-checking with NLP techniques. However, their research points out more of the fact that automated fact-checking is less reliable when compared to manual fact-checking. Das et al. (2023)'s solution for this limitation is to develop a hybrid system for automatic fact-checking that would use humans in the process alongside computers. However, in the long run, especially when we consider the possibility of employing large language models to generate false news, using humans in the detection process, although likely still achieving some decent results, would most likely fail to perform as effectively as is required.

Waikhom [10] used ensemble machine learning methods, such as XGBoost, bagging, random forest (RF), extra trees, and gradient boost. The methods they used allowed them to achieve relatively high accuracy scores in classification when using the LIAR dataset [4]. Ahmad [27] also used ensemble methods with machine learning in their research about false news detection. They also achieved very good results with ensemble methods. Waikhom [10] and Ahmad's [27] research and results conclude that machine learning algorithms work well for false news detection when implemented with ensemble learners. However, Waikhom [10] and Ahmad's [27] studies might overemphasize the general effectiveness of ensemble learners while ignoring potential drawbacks, such as higher processing requirements, the risk of over fitting, and the challenge of sustaining performance over more diverse datasets.

Gundapu [12] researched false news detection for COVID-19-related news. They used classic ML models, deep learning models, and transformer models to conduct their research. They achieved the best results when they developed an ensemble model consisting of all three different transformer models that they used (BERT, ALBERT, and XLNET). With this ensemble model, they were able to receive an accuracy score of 98%. However, the usage of one specific dataset may have contributed to the model's high accuracy, potentially limiting its usefulness in broader scenarios.

Wu [28] conducted their research on multimodal (text and image) false news detection. They proposed a multimodal co-attention network-based model to better include text and images together for false news detection. Their model first extracted visual features from images and then textual features from the text, fusing these extracted features that then can be used to detect false news. Their model was able to achieve good accuracy results on the two datasets that they used for their research. On the first dataset that they used, they

achieved an accuracy score of 80%, and on the second, they achieved an accuracy score of 89%. Nadeem [29] recently concluded research on utilizing visual features for false news detection as well. They proposed a multimodal extreme fake news detection (EFND) that gathers context, social context, and visual data to create a multimodal vector. The results they achieved were high, with the accuracy score being 98% and 99% on different datasets. Wu's [28] model which combined text and images achieved good results, but Nadeem's [29] model, which included additional social and contextual data, achieved higher accuracy results. This indicates that multimodal approaches—while still efficient—benefit from additional context to combat their complexity that can impact their ability to generalize and scale effectively.

2.1.2. Linguistic and Textual Analysis of False News

Singh [30] used linguistic analysis alongside ML in their research. Their research provides interesting information about linguistic differences between false and real news. From their research, one can learn that in general, false news tends to be shorter, show less expertise or confidence, appear negative in tone, and show less analytical thinking. However, their research shows that the package they used for linguistic analysis, linguistic analysis and word count (LIWC), associates the language found in false news with higher authenticity. In LIWC, a higher authenticity score is assigned when the language is more personal and disclosing. In comparison, a lower authenticity score is assigned when the language is more guarded and distanced. This could explain why people can be tricked into believing false news, or it could reveal a potential misalignment in LIWC, where personal language can misleadingly enhance credibility.

Ahmed [31] proposed a detection model that combines text analysis using n-gram features and term frequency with ML classification. In addition, they introduce a novel n-gram model in their research, which generates distinct sets of n-gram frequency profiles from their trained data to distinguish between false and true content. Their research revealed that linear-function-based classifiers outperformed non-linear classifiers. The researchers also discovered that increasing the n-gram size affected detection accuracy. This suggests that the language used in false news is inconsistent.

2.1.3. Automatically Generated Text Detection

Mitrović [32] researched detecting short texts generated by ChatGPT using a transformer-based model. In their research, they also analyzed the language generated by ChatGPT and concluded that "ChatGPT's writing is polite, without specific details, using fancy and atypical vocabulary, impersonal, and typically it does not express feelings." (e.g., [32], P. 1). The research focused on restaurant reviews generated by ChatGPT, and the goal was to classify the reviews according to whether they were created by a human or ChatGPT. The research achieved a good 79% accuracy, even though the research stated that the transformer-based model had problems with differentiating between human and ChatGPT-generated reviews.

2.1.4. False News Detection Based on User Interaction

Tacchini [33] proposed an idea regarding hoax and false news detection, where the nature of a Facebook post could be determined by the users who "like" the posts. The baseline for their research was that a user who "likes" a post determined as a hoax, is anticipated to "like" even more hoax posts. They would analyze a post according to the users who "liked" the post, and if there was a high enough amount of users who had previously "liked" several posts determined as hoaxes, the current post being analyzed would be determined to be a hoax as well. For their experiments, they used two different classification techniques: logistic regression and boolean label crowdsourcing (BLC). Both of their techniques achieved very high results in detecting whether a Facebook post could be determined as a hoax or a non-hoax, suggesting that analyzing the users interacting with Facebook posts can accurately determine the nature of the Facebook posts. While

Tacchini's [33] approach effectively employs user interaction to detect hoax posts with high accuracy, its reliance on user behavior, which could be influenced by various factors unrelated to the post's authenticity, may limit its robustness in reality.

Del Tredici [34] used linguistic analysis and user detection to detect false news. They proposed a model that would generate representations of social media users based on the language they use and the news they spread, which would then be used to detect false news. The model was created by using convolutional neural networks (CNNs), which are ideally suited for text classification. In their study, they analyzed the language commonly used by those who share false news. The study concluded that the language used by users who spread false news is consistent, which, in turn, makes it easier to detect news based on the people who share them, like in Tacchini's [33] research. However, Del Tredici's [34] approach may struggle to account for the variety and developing tactics of misinformation across different social media contexts.

2.2. Datasets

In this sub-section, we go through the datasets used in this article. These datasets include three datasets created by the previous research, namely, LIAR [4], FakeNewsNet [35], and Twitter15 [36], as well as a novel dataset that we built by using ChatGPT. We analyze the sizes of the datasets, the data elements used in the datasets, and the different values of the data elements. We have included a summary of all the datasets that can be seen in Table 4.

Table 4. A summary of all used datasets.

Datasets	All Samples	True Samples	False Samples	Information Type
LIAR	12,851	7134	5707	News related to politics.
FakeNewsNet	23,921	6480	17,441	News related to politics and celebrity gossip.
Twitter15	1490	372	370	Rumors spread on Twitter.
Novel ChatGPT	300	100	200	Automatically generated false and real news articles.

2.2.1. Liar Dataset

The LIAR dataset [4] is a benchmark dataset created for false news detection. It contains 12.8 thousand real-world manually labelled short statements that were collected from PolitiFact.com with various contexts. The dates for the statements are primarily from 2007–2016. The dataset also includes an analysis report and links to source documents for each statement, as well as information about the speaker, the speaker's job title, subject, political party affiliation, the credit history of the speaker, and the context for each statement.

The dataset has six different labels to determine the truthfulness ratings. These labels are as follows: pants-fire, false, barely-true, half-true, mostly-true, and true. The pants-fire label represents a completely false statement, and the true label represents a completely true statement. The distribution of cases for each label is balanced, except for pants-fire which has significantly fewer cases compared to other labels. The pants-fire label has 1050 cases whereas the other labels have cases ranging from 2063 to 2638 cases.

The statements have 732 different subject types ranging from various topics, with the most frequent subject being healthcare, and it appears in the dataset 5 times. The average statement is 17.9 tokens long. Most of the speakers of the statements are U.S. politicians, but other speaker types are also included such as journalists, social media users, and Internet newspapers. Overall, there are 2910 unique speakers, and each speaker appears 3.5 times on average in the dataset. The most common speaker in the dataset is Barack Obama, who appears in the dataset 5 times. Table 5 shows an example of a statement from the LIAR dataset [4].

Table 5. An example of a randomly chosen statement from the LIAR dataset. The label history shows how many times the speaker has made a statement that belongs to one of the six different label cases in the dataset.

Elements	Value of Elements
ID	8303
Label	half true
Statement	Tuition at Rutgers has increased 10 percent since Gov. Chris Christie took office because he cut funding for higher education.
Subject	education, state finances
Speaker	Barbara Buono
Job title	State Senator
Party affiliation	democrat
Label history	3, 1, 4, 4, 1
Context	a speech to students at the Rutgers New Brunswick campus

2.2.2. FakeNewsNet Data Repository

The FakeNewsNet [35] is a data repository that contains two datasets. The datasets were collected from PolitiFact.com and GossipCop. The datasets include the collected news articles, social context, and information about users who interacted with the article on social media. The inclusion of user interaction information makes this data repository useful when detecting false news from social media. As shown by Tacchini [33] and Del Tredici [34] the inclusion of user analysis is an excellent way to detect false news that is spread in social media. The datasets contain source URLs to the news articles, the title of the news, and the tweet IDs of users who interacted with the article on Twitter.

The datasets are different in size, where the dataset collected from GossipCop is considerably larger than the dataset collected from PolitiFact.com. The distribution of false and true news in the datasets is imbalanced, especially in the dataset collected from GossipCop. The dataset collected from PolitiFact.com contains 432 news articles labeled as false and 624 news articles labeled as real. The dataset collected from GossipCop contains 6048 news articles labeled as false and 16,817 labeled as real.

The articles in the PolitiFact.com dataset focus on political issues, whereas the articles in the GossipCop dataset contain news about celebrities. GossipCop used to be a website to fact-check articles and stories related to the entertainment industry. As GossipCop mainly focused on false stories, it provides a reason why the imbalance between real and false articles is large in the dataset.

The average title length in the PolitiFact.com dataset is 10.74 tokens, and in the GossipCop dataset, it is 10.067 tokens. In the PolitiFact.com dataset, a news article was interacted with by 1.329 different users on average, and in the GossipCop dataset, the same average was 1.064. Table 6 shows an example of statements found in the FakeNewsNet [35] data repository.

2.2.3. Twitter15

Twitter15 [36] includes 1490 Twitter stories posted until March 2015. The stories were collected from Snopes.com and Emergent.info. The dataset is used for rumor detection on Twitter posts. The distribution of false and true events is similar in size. The dataset contains 372 events determined as true rumors, 370 events determined as false rumors, 374 events determined as non-rumors, and 374 events that could not be verified. The labels used to differentiate the stories are non-rumor, true, false, and unverified.

Table 6. Examples of randomly chosen false and real news articles from both FakeNewsNet’s datasets.

Elements	Value of Elements
ID	politifact182
News URL	http://www.gao.gov/new.items/d071195.pdf (accessed on 12 August 2024).
Title	US Government Accountability Office Report to Congressional Committees
Tweet IDs	956894522511736832
Label	real
ID	politifact14944
News URL	http://thehill.com/homenews/senate/369928-who-is-affected-by-the-government-shutdown (accessed on 12 August 2024)
Title	Who is affected by the government shutdown?
Tweet IDs	954602090462146560 954602093171609600 954650329668349954
Label	false
ID	gossipcop-897603
News URL	https://www.teenvogue.com/story/selena-gomez-not-changing-blonde-hair (accessed on 12 August 2024)
Title	Selena Gomez Is Going To Keep Her Blonde Hair
Tweet IDs	936830208857878528
Label	real

As the dataset contains stories posted on Twitter, the genre of the stories varies a lot. The most common words found in the dataset were Paul, shot, new, police, says, killed, war, Ferguson, died, and Obama. The average text length of a post was 10.2 tokens. Table 7 shows an example of a statement found in the Twitter15 dataset [36].

Table 7. An example of a randomly chosen news article from the Twitter15 dataset.

Elements	Value of Elements
ID	693560600471863296
Events	miami was desperate for a turnover. instead, nc state got this dunk. and a big upset win: URL
Veracity	non-rumor

2.2.4. Novel ChatGPT-Generated Dataset

For experimental purposes, We created a novel dataset that includes news articles generated by ChatGPT to test how well classifiers would handle artificially generated text. The dataset consists of an ID, title, text, author/source, and the labels false and true. Additionally, the texts labeled as true are provided with a link to the original article. We created 200 false news articles for the dataset and 100 true articles.

Creating false news articles with ChatGPT was a very specific and unusual task. By default, ChatGPT refuses the creation of false news articles. This, however, can be easily manipulated with wordplay and tricking ChatGPT into writing articles that contain misleading information. In Figure 3 we can observe our initial failure in creating false news with ChatGPT, and in Figure 1, we can observe how we succeeded in manipulating ChatGPT to generate false news for us. The text generated by ChatGPT is believable, and creative, and does not contain language usually found in false news, such as exaggerated language. The example of how to generate false news articles with ChatGPT shows how easily it can be used to generate untruthful information in general.

The so-called truthful news for this dataset was created by first summarizing actual news articles from credible news sources, mainly from Reuters (<https://www.reuters.com/>, accessed on 12 August 2024), Helsinki Times (<https://www.helsinkitimes.fi/>, accessed on 12 August 2024), The Kyiv Independent (<https://kyivindependent.com/>, accessed on 12 August 2024), and NHK WORLD-JAPAN (<https://www3.nhk.or.jp/nhkworld/>, accessed on 12 August 2024), then making ChatGPT write the news article again using the summary that was created. We also made sure that each generated article was truthful towards its

original context, and that it did not contain false information made up by ChatGPT. We also provide a link to the original article that was summarized and rewritten by ChatGPT. Table 8 shows an example of a statement found from the novel ChatGPT-generated dataset.



Figure 3. An example of failure when trying to trick ChatGPT into writing a false news article.

We discovered that the articles in other datasets are not particularly diverse. Many of them are primarily focused with the United States and its political climate, leaving other parts of the world and different topics virtually unaddressed. For this issue, we included a wide range of topics, including sports, economics, medicine, and crimes from several countries in this dataset. This adjustment was made to ensure more diversity in the data, preventing the trained model from simply over fitting to a certain politics-related term in classification.

Table 8. An example of an entry in the novel ChatGPT-generated dataset.

Elements	Value of Elements
ID	0
title	Fed plans broad revamp of bank oversight after SVB failure
text	The Federal Reserve could make a significant impact on its supervisory practices by rapidly implementing mitigants in response to serious issues regarding capital, liquidity, or management, according to a senior Fed official. . .
source	Reuters
label	true
original article	https://www.reuters.com/business/finance/fed-plans-broad-revamp-bank-oversight-after-svb-failure-2023-04-28 (accessed on 12 August 2024)

2.3. Applied Methods

This sub-section goes through the different feature extraction and classifying methods used in the research.

2.3.1. Feature Extraction

CountVectorizer turns given textual data into a vector based on the frequency (count) of each word in the text. The created vector is represented as a sparse matrix, where each of the words is stored using index values determined by alphabetical order. Using CountVectorizer makes it easy to use textual data directly in ML models in text classification tasks.

TF-IDF (term frequency-inverse document frequency) is a technique used to determine the importance of words found in used documents. The TF-IDF score of a word is determined according to how many times the word appears in the document. A word that appears less frequently obtains a higher score than a word that appears more frequently. This is because a word that appears less frequently in a sentence is seen as more important since it usually gives a better context to the sentence.

The TF-IDF score is a combination of two calculations: the term frequency (TF) and the inverse document frequency (IDF). The TF score is calculated by dividing the number

of occurrences of a word in a document by the total number of words in that document. The IDF score is calculated by dividing the total number of documents by the documents containing a certain word. The TF and IDF scores are then multiplied, and finally, the TF-IDF result can be obtained.

In this research, we applied the TF-IDF technique using the `TfidfTransformer`. It is used to transform a count matrix into a matrix of TF-IDF scores. `TfidfTransformer` takes the count matrix as an input and applies the TF-IDF technique to convert the matrix into a weighted representation.

With `CountVectorizer` and `TfidfTransformer`, we created a bag-of-words model. Bag-of-words is a common technique in NLP, where the used textual data are turned into numerical features that ML algorithms are capable of processing.

2.3.2. Classifying Methods

Random forest (RF) is an ensemble learning method that uses a combination of multiple decision trees when training. When used in a classification task, RF generates an ensemble of trees that predict the classification result by casting a vote to determine the most popular class. In RF, each tree in the ensemble is built from a sample drawn with a replacement from a training set—which, in the case of our experiments, are the datasets that we used [37]. The research conducted by Waikhom, et al. [10] and Ahmad, et al. [27], which we discussed more in detail in Section 2.1, showed that ensemble learning methods achieved high accuracy scores for false news detection. This provides a reason why we chose to include RF as one of the methods for our classifying experiments.

Naive Bayes (NB) methods are a set of supervised learning algorithms based on applying Bayes' theorem with strong (naive) independence assumptions between features given the value of a class variable. In this experiment, we used the `MultinomialNB` algorithm in our classification experiments. `MultinomialNB` predicts the probability of general labels for a given text, applying those features in a multinomial function based on Bayes' theorem. NB algorithms are often used as a baseline for text classification tasks; they work well with smaller datasets and are faster to train when compared with other popular classifiers [38].

Logistic Regression (LR) is a statistical model that calculates the probability of a binary outcome, based on prior observations from a dataset. LR can consider multiple input criteria for the eventual outcome of a prediction. In the instance of false news detection, the LR model could be capable, for example, of taking into consideration the history of label distribution (true/false) of news articles written by a reporter. Based on the historical data, the LR model calculates the score for a new case based on its probability of receiving either one of the labels [39]. LR has been successful for false news detection purposes previously [31,33,40].

Support vector machine (SVM) is a machine learning algorithm that uses supervised learning models to analyze data for classification, regression, and outlier detection. SVM works by creating a hyperplane that separates the used training data into classes, e.g., false/true. The new data are then fitted into the same space as the old data, and the class prediction for the new data is determined by which side of the previously created hyperplane they fall [41]. SVM has been used widely in false news detection, often being able to receive high accuracy results [23,42,43].

k-Nearest Neighbors (kNN) is a non-parametric supervised learning method used for classification and regression tasks. The kNN classifier takes the k closest training data as an input, and the input data are then classified based on a plurality vote of their neighbors, meaning that an object currently being classified is assigned to a class according to its k nearest neighbors, where k is a positive integer [44]. In our research, we used $k = 3$, where an object is classified according to what is the most common class among its neighbors. Like NB, kNN is often used as a baseline for text classification tasks and is fast and simple to train [38].

Multi-layer Perceptron (MLP) is a feed-forward artificial neural network, typically consisting of three different interconnected layers: input layer, hidden layer, and output layer. In MLP classification tasks, the input data are passed through the layers, where each layer solves a specific part of the task. The output of the solved result is passed through the layers until the result is determined, e.g., whether an article is false or not [45]. For our research, we implemented our MLP classifier with the suggested settings given by Scikit-learn's user guide for neural networks (https://scikit-learn.org/stable/modules/neural_networks_supervised.html#classification, accessed on 12 August 2024). MLP has shown good results in the previous research about false news detection [11,46], and it has been useful for more complex tasks than simple text-based classification [29].

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression tasks. They have a hierarchical, tree structure that consists of a root node, branches, internal nodes, and leaf nodes. The root node and the internal nodes represent the base problem given by a classification task, e.g., is this article false or true; branches represent the outcome of the problem; and the leaf nodes represent the final decision after calculating all the possible outcomes, e.g., this article is false [47]. DTs are simple and easy to understand, and because of their structure, they are capable of handling multi-output classification problems. For false news detection, DTs have been capable of providing good accuracy results [46,48].

Boosting in ML refers to a set of ensemble algorithms, such as adaptive boosting (AdaBoost), gradient boosting, or extreme gradient boosting (XGBoost). For this research, we implemented the AdaBoost algorithm. AdaBoost works by training a classifier on a dataset, and the classifier is given a weight according to the performance. AdaBoost gives a higher weight to items that have been classified incorrectly so that the incorrect classification can be corrected. This process is then repeated until the actual values and predicted values reach an acceptable threshold [49]. The AdaBoost algorithm is simple to use, yet it can achieve good accuracy results for false news detection tasks [10].

Stochastic Gradient Descent (SGD) is an algorithm used to minimize the loss functions of a linear classifier model. It is often used in large-scale machine learning problems usually encountered in text classification and natural language processing. SGD is an optimization method, used to train a model to find the optimal set of parameters for the model. In this research, the model we trained using SGD was the "Modified Huber" loss function. SGD is efficient and easy to implement, and it works well with large datasets [50]. SGD has been capable of achieving good accuracy results for false-news-detection-related tasks [48].

Passive Aggressive (PA) is an algorithm that is part of a group called linear models, where the target value is expected to be a linear combination of the features [51]. More specifically, PA is used with binary classification tasks, usually when the used data are potentially noisy or they might change over time [52]. Passive aggressive algorithms have been widely used for false-news-detection-related tasks, often with good results [53,54].

3. Results

In this section, we introduce the experiments conducted on the four datasets that were introduced in Section 2.2 and their results.

3.1. Experiment 1: Twitter15 Dataset

For our first set of experiments, we used the Twitter15 dataset [36]. This dataset was already divided into separate "test" and "train" files. As this dataset was created for rumor detection, we used this dataset to compare the differences and similarities between false news and rumors. As presented in Table 1, the definitions of rumors and false news have a slight overlap; we thought that it would be interesting to conduct a classifying experiment using both rumors and false news. In this experiment, we first conducted a classifying experiment only using the Twitter15 dataset [36]. Then, we experimented using the LIAR dataset [4] as train data and the Twitter15 dataset [36] as test data. For the final experiment, we used the Twitter15 dataset [36] as train data and the LIAR dataset [4] as test data.

We first built a bag-of-words model with *CountVectorizer* and *TfidfTransformer* to use the ML classifiers. For this experiment, to obtain the best comparisons between all the different methods, we used all the methods introduced in Section 2.3.2. For the later experiments, we chose the three best-performing methods from this experiment.

When only the Twitter15 dataset [36] was used in the tests, the split ratio was around 75:25, as specified by the original authors. Using the LIAR dataset [4] as either the “train” or “test” file, along with the Twitter15 dataset [36], the split ratio was 90:10.

Our first experiment was first conducted only using the Twitter15 dataset [36] on the classifiers. This allowed us to create a decent baseline for our classifiers, which we could then compare to the classifying results obtained later by more complex studies. The results of this experiment can be seen in Table 9, and the ROC curve for this experiment can be observed in Figure 4. From Table 9, we can observe that all the classifiers performed well on average, as all the classifiers could achieve accuracy results of over 50%. The worst-performing classifiers were AdaBoost and MLP, which achieved accuracy scores of 54% and 68%. These classifiers did not perform well with AUC, FPR, and FNR scores either. Especially when looking at the FPR and FNR scores, we can observe that these classifiers had the most tendencies to misclassify the labels. The best-performing classifiers were the passive aggressive classifier and SVM, which achieved accuracy scores of 87%.

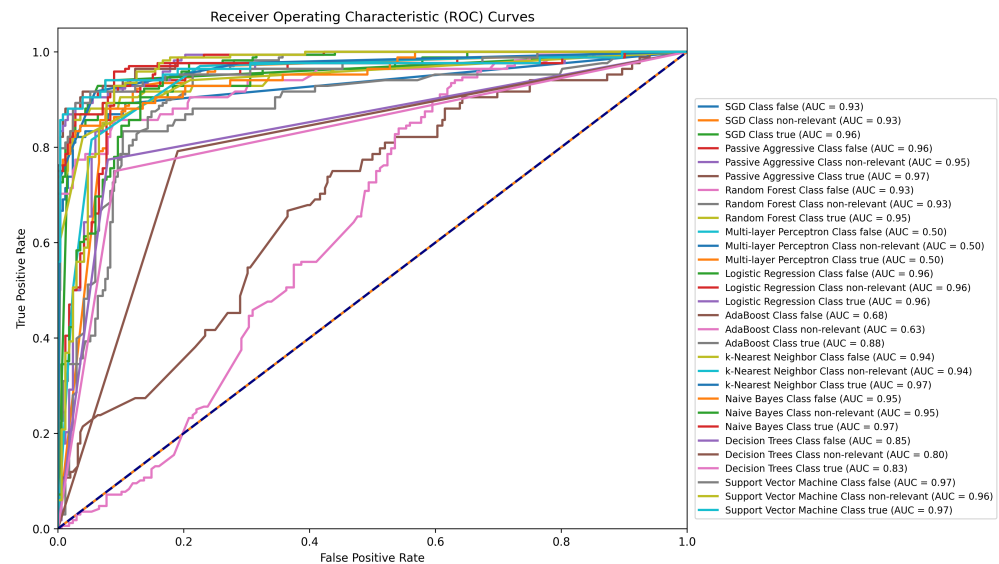


Figure 4. ROC curve for the classifying results conducted only using the Twitter15 dataset.

Table 9. Performance report on the classifying results conducted only using the Twitter15 dataset.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
SGD	0.84	0.84	0.84	0.84	0.93	0.04	0.15
PA	0.86	0.87	0.87	0.87	0.95	0.04	0.12
RF	0.81	0.82	0.80	0.80	0.94	0.06	0.43
MLP	0.68	0.70	0.68	0.68	0.50	0.33	0.66
LR	0.85	0.86	0.85	0.86	0.96	0.03	0.31
ADA	0.55	0.54	0.54	0.53	0.72	0.00	0.97
kNN	0.79	0.79	0.79	0.78	0.95	0.05	0.15
NB	0.80	0.80	0.80	0.80	0.96	0.06	0.36
DT	0.71	0.72	0.72	0.72	0.83	0.11	0.22
SVM	0.87	0.87	0.87	0.87	0.96	0.03	0.16

The second and third experiments were conducted using the Twitter15 dataset as train data and the LIAR dataset as test data, and vice versa. The results for these experiments can be seen in Tables 10 and 11, and the ROC curves for these experiments can be observed in

Figures 5 and 6. These tables show that testing and training with various false information datasets did not produce good results. This is understandable given that the vocabulary used in false news differs from that used in rumors.

Table 10. Performance report on the classifying results conducted using Twitter15 as train data and LIAR as test data.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
SGD	0.34	0.32	0.32	0.28	0.49	0.33	0.67
PA	0.35	0.35	0.34	0.29	0.49	0.32	0.69
RF	0.36	0.31	0.33	0.22	0.50	0.27	0.73
MLP	0.38	0.13	0.33	0.18	0.50	0.33	0.66
LR	0.37	0.28	0.33	0.21	0.51	0.27	0.74
ADA	0.36	0.32	0.33	0.24	0.50	0.00	1.00
kNN	0.35	0.34	0.34	0.32	0.52	0.28	0.70
NB	0.37	0.18	0.33	0.18	0.50	0.30	0.62
DT	0.36	0.32	0.32	0.27	0.49	0.33	0.67
SVM	0.36	0.31	0.34	0.26	0.50	0.32	0.68

Table 11. Performance report on the classifying results conducted using Twitter15 as test data and LIAR as train data.

Method	Acc.	Precision	Recall	F1	AUC	FPR	FNR
SGD	0.26	0.28	0.27	0.27	0.41	0.30	0.78
PA	0.24	0.24	0.24	0.23	0.40	0.49	0.67
RF	0.39	0.34	0.33	0.32	0.49	0.01	0.94
MLP	0.50	0.17	0.33	0.22	0.50	0.00	1.00
LR	0.28	0.28	0.28	0.27	0.43	0.05	0.97
ADA	0.40	0.30	0.31	0.29	0.44	0.00	1.00
kNN	0.31	0.34	0.32	0.31	0.51	0.25	0.75
NB	0.33	0.29	0.31	0.29	0.44	0.02	0.99
DT	0.35	0.31	0.31	0.21	0.48	0.33	0.69
SVM	0.27	0.29	0.28	0.27	0.42	0.44	0.66

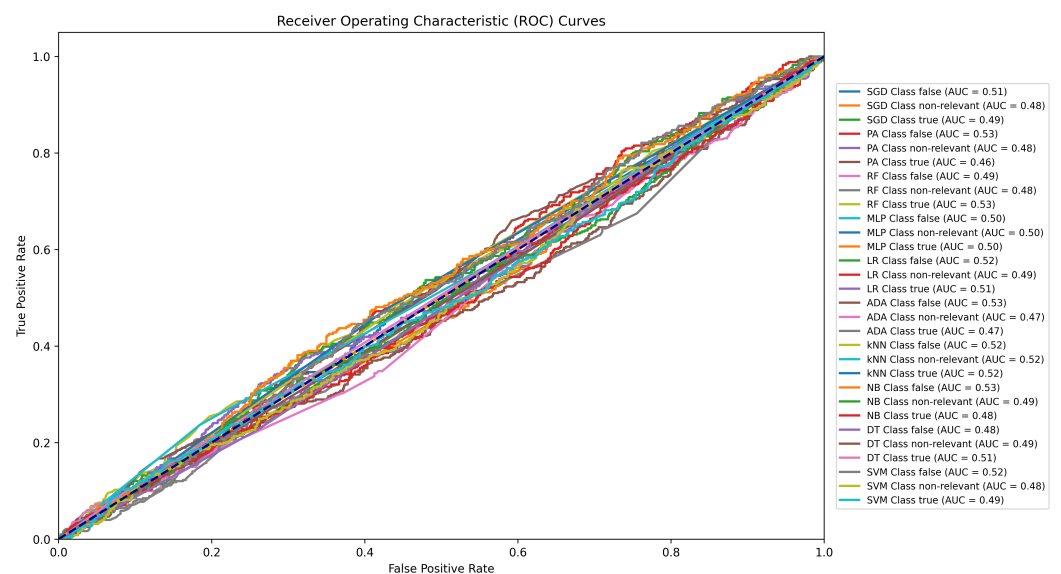


Figure 5. ROC curve for the classifying results conducted using the Twitter15 dataset as train data and the LIAR as test data.

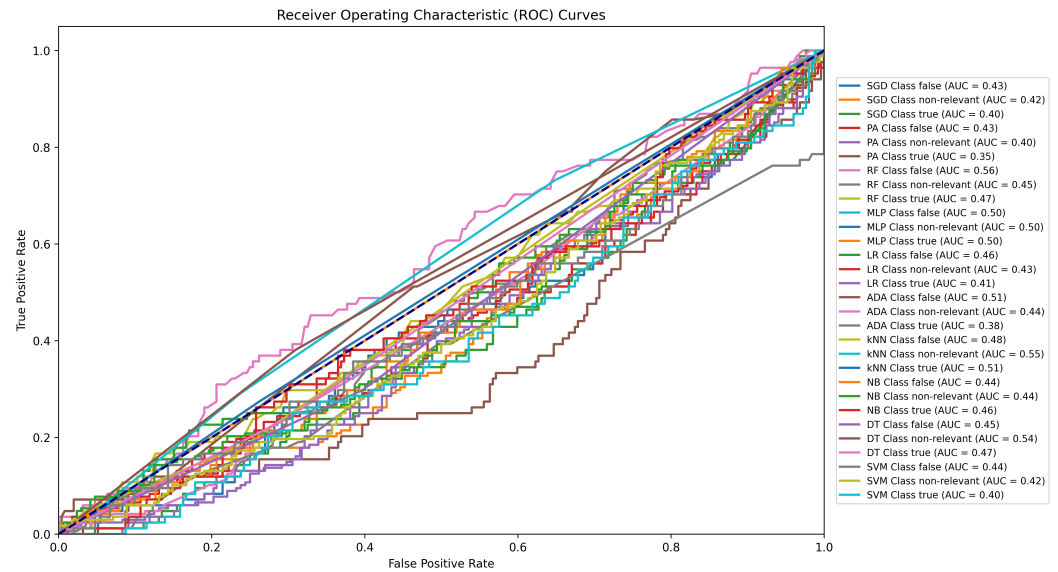


Figure 6. ROC curve for the classifying results conducted using the Twitter15 dataset as test data and the LIAR as train data.

Based on these classification experiments, we concluded that RF, PA, and SVM were the best-performing methods and the most consistent with their precision, recall, F1, AUC, FPR, and FNR scores. We use these three methods for the later classification experiments for these reasons.

3.2. Experiment 2: LIAR Dataset

We used the LIAR dataset [4] for our second set of experiments. This dataset was already divided into training, validation, and testing files. We combined the validation and training data for these experiments, as the inclusion of the validation file was unnecessary for our experiments.

To be able to use the ML classifiers, we first built a bag-of-words model with CountVectorizer and TfidfTransformer, like in the previous experiment. For the experiments, we used three different ML classification algorithms: *passive aggressive classifier* (PA), *random forest* (RF), and *support vector machine* (SVM).

As the LIAR dataset originally separates false and true news into six different labels, we concluded two different types of classification experiments in this dataset. At first, we concluded a six-label classification experiment, where we classified articles with the six different labels originally provided by the dataset’s author. After that, we concluded a series of different binary classification experiments, where we reduced the labels of the dataset into only “true” and “false”. We reduced the labels in five different ways: all labels except “true” are labeled as “false”; labels “pants-fire” and “false” are labeled as “false” and the rest as “true”; all labels except “pants-fire” are labeled as “true”; labels are split from the middle into “true” and “false”; and labels “true” and “mostly-true” are labeled as “true” and the rest as “false”. We used these labeling methods to determine when non-binary classification works best for the LIAR dataset and which labels make the language the most differentiable.

When performing the six-label classification, a split ratio of 90:10 was employed as that is how the dataset was split by the original author. During the multiple binary-label studies, we divided the dataset into a split ratio of 80:20.

The results for the six-label classification experiment can be found in Table 12, and the ROC curve for this experiment can be observed in Figure 7. From these results, we can observe that all the used methods achieved low scores overall. The RF method achieved the best accuracy score; however, the FPR and FNR scores for this method were not good and showcased severe bias for detecting negative classes only. The PA and SVM

methods performed more consistently, especially in terms of FPR and FNR scores, but their performance was not high either.

After the six-way classification experiment, we concluded a series of binary classification experiments. The results for the binary classification experiments can be found in Tables 13–17. The ROC curves for these experiments can be found in Figures 8–12. From these tables, we can observe that the classifiers were most efficient when only the label “pants-fire” was considered as false and least efficient when the labels “mostly-true” and “true” were considered true and the rest were false. Overall, in the binary classification experiments, the classifiers were more accurate than in the six-label experiment. The methods in general performed equally, but in some binary classification experiments, the RF method again showed severe bias towards detecting only negative classes, whereas PA and SVM were more consistent with their FPR and FNR scores.

Table 12. Performance report for the six-label classification experiment.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.21	0.20	0.20	0.20	0.56	0.34	0.58
RF	0.25	0.25	0.22	0.21	0.56	0.00	1.00
SVM	0.23	0.22	0.22	0.21	0.57	0.28	0.64

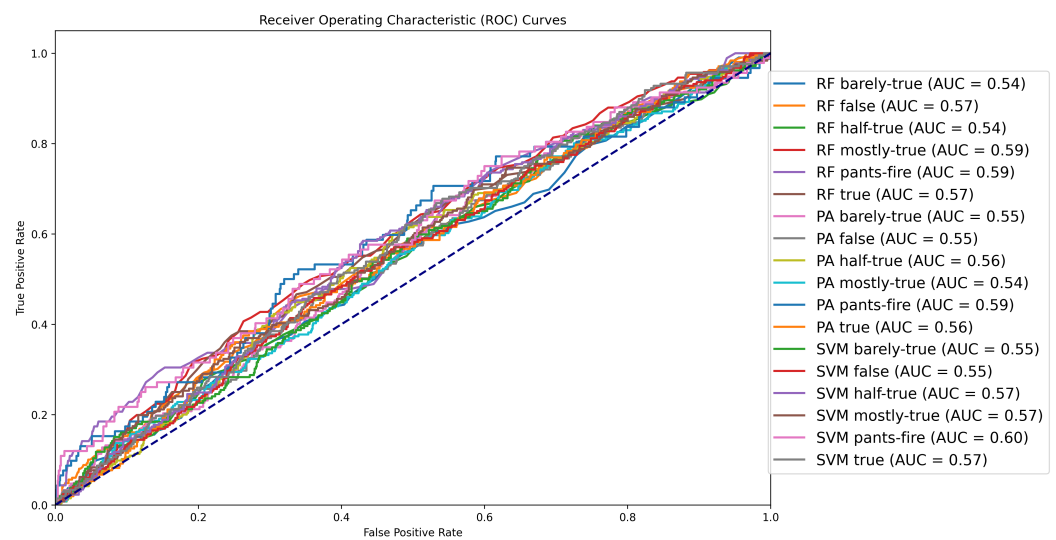


Figure 7. ROC curve for the six-label classification experiment.

Table 13. Performance report for the binary classification experiment where all the labels except “true” are re-labeled as false.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.76	0.52	0.51	0.51	0.57	0.39	0.46
RF	0.83	0.91	0.50	0.46	0.58	0.00	0.99
SVM	0.81	0.51	0.50	0.47	0.58	0.25	0.67

Table 14. Performance report for the binary classification experiment where labels “pants-fire” and “false” are re-labeled as false and the rest as true.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.79	0.51	0.51	0.51	0.54	0.48	0.44
RF	0.86	0.43	0.50	0.46	0.50	0.99	0.01
SVM	0.85	0.51	0.50	0.48	0.54	0.54	0.39

Table 15. Performance report for the binary classification experiment where only the label “pants-fire” is considered false, and the rest are true.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.88	0.53	0.51	0.51	0.53	0.50	0.46
RF	0.91	0.45	0.50	0.48	0.57	1.00	0.00
SVM	0.91	0.45	0.50	0.48	0.53	0.65	0.29

Table 16. Performance report for the binary classification experiment where labels “barely-true”, “pants-fire” and “false” are considered false, and labels “half-true”, “mostly-true” and “true” are considered true.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.56	0.55	0.55	0.55	0.59	0.37	0.53
RF	0.59	0.59	0.59	0.59	0.63	0.45	0.35
SVM	0.58	0.58	0.58	0.58	0.60	0.25	0.62

Table 17. Performance report for the binary classification experiment where labels “mostly-true” and “true” are considered true and the rest false.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.56	0.49	0.49	0.49	0.49	0.46	0.53
RF	0.65	0.50	0.50	0.44	0.53	0.07	0.93
SVM	0.60	0.51	0.50	0.50	0.50	0.40	0.59

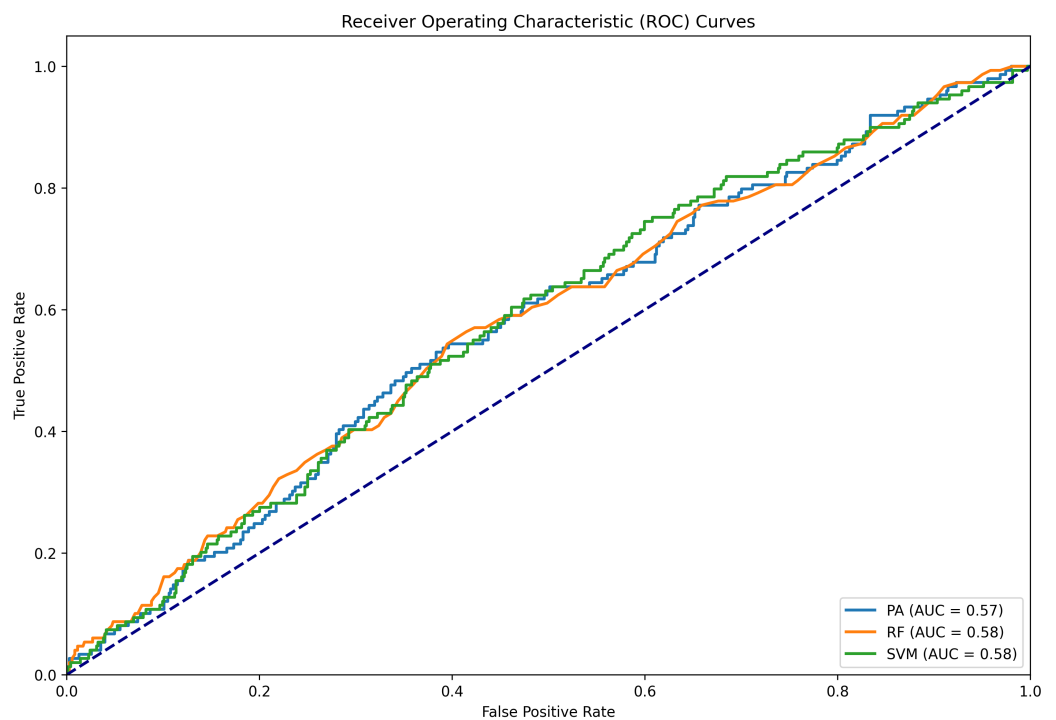


Figure 8. ROC curve for the binary classification experiment where all the labels except “true” are re-labeled as false.

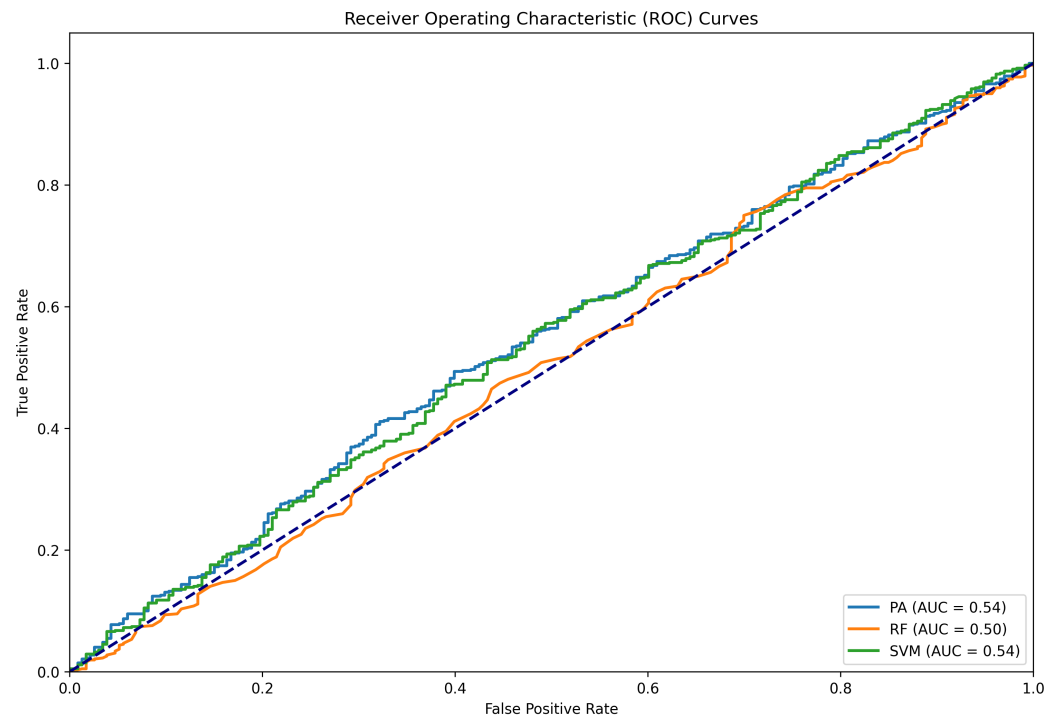


Figure 9. ROC curve for the binary classification experiment where labels “pants-fire” and “false” are re-labeled as false and the rest as true.

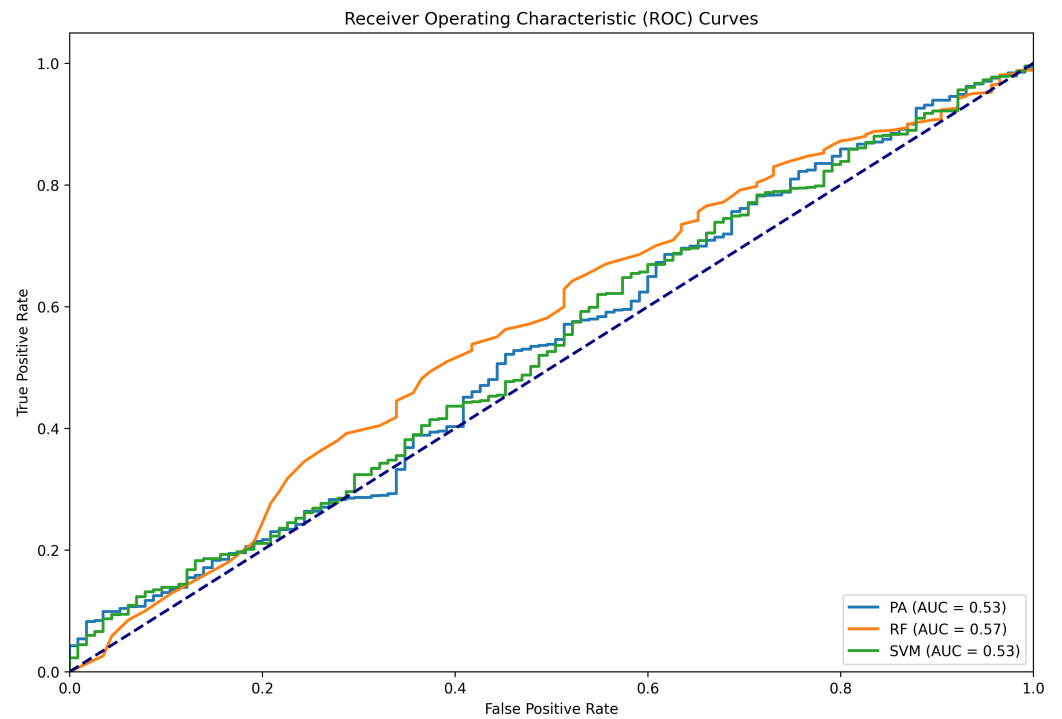


Figure 10. ROC curve for the binary classification experiment where only the label “pants-fire” is considered false, and the rest are true.

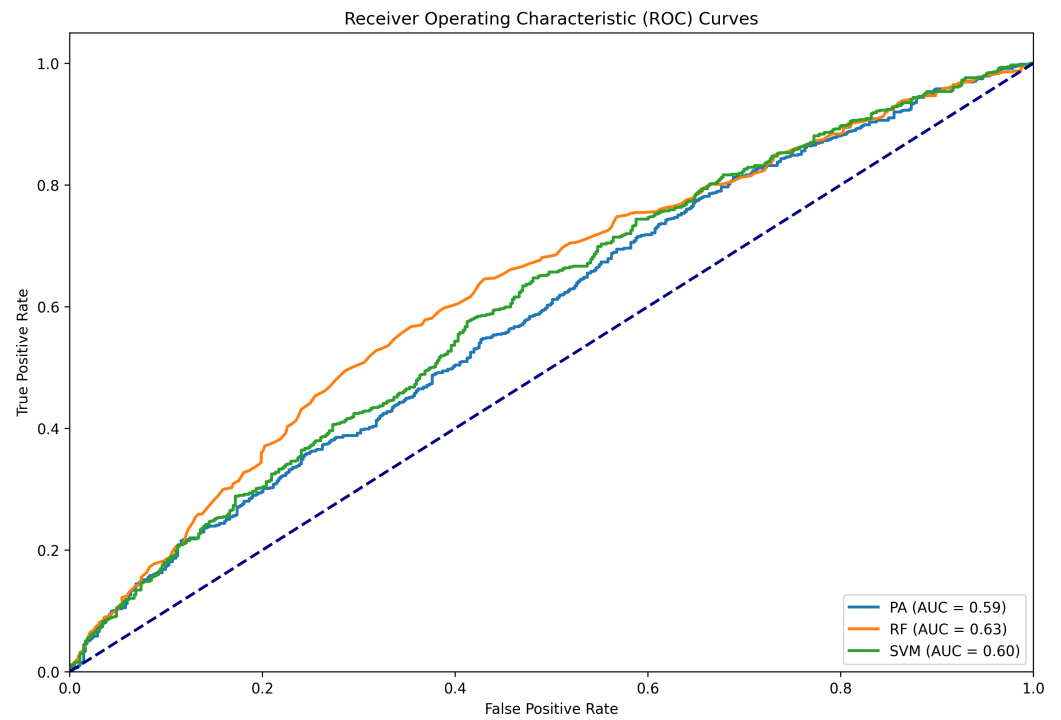


Figure 11. ROC curve for the binary classification experiment where labels “barely-true”, “pants-fire” and “false” are considered false, and labels “half-true”, “mostly-true” and “true” are considered true.

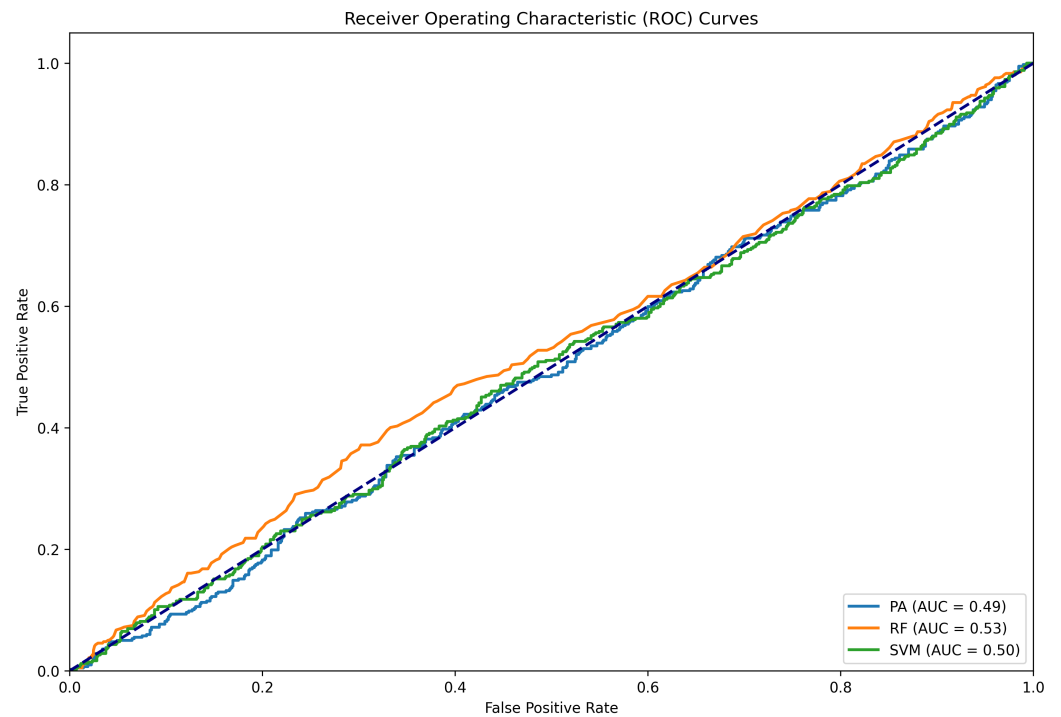


Figure 12. ROC curve for the binary classification experiment where labels “mostly-true” and “true” are considered true and the rest false.

3.3. Experiment 3: FakeNewsNet Data Repository

Our third set of experiments was conducted on the FakeNewsNet data repository [35]. As explained in Section 2.2, this data repository consists of two datasets, one that has a focus on political articles and one that has a focus on celebrity gossip. The datasets were split into “real” and “false” files. Before the experiments, we added the labels “true” and

“false” to all the data files. After that, we combined the “real” and “false” files according to the source where the dataset was collected from.

Like in the previous experiments, we first built a bag-of-words model using CountVectorizer and TfidfTransformer. The methods we used for these experiments were the same as in Section 3.2.

We concluded four different experiments with this data repository, one where we would use the PolitiFact dataset as the “train” file and the GossipCop dataset as the “test” file, and vice versa, as well as an experiment where we split the datasets into “train” and “test”. As both datasets are very different from one another in terms of content, we wanted to experiment with how using the two datasets together would affect the possible results.

When the two datasets were utilized as either the “train” or “test” files, the split ratio in these studies was approximately 90:10. Because of the differences in size between both datasets, it was difficult to strike a good balance between maintaining the “train” and “test” as similar in both studies.

In the studies where the datasets were not utilized to train and test each other but rather were separated independently into “train” and “test”, the split ratio was 80:20.

The results for all the classification experiments can be seen in Tables 18–21. The ROC curves for the results of the classification experiments can be observed in Figures 13–16. From these tables, we can observe that the classification methods were the least efficient when training with the PolitiFact dataset and testing with the GossipCop dataset. The highest accuracy scores were achieved when splitting the GossipCop dataset into “train” and “test”.

Table 18. Performance report when the PolitiFact dataset was used as “train” and the GossipCop dataset was used as “test”.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.42	0.50	0.49	0.41	0.50	0.63	0.37
RF	0.64	0.51	0.51	0.51	0.53	0.74	0.26
SVM	0.44	0.50	0.50	0.43	0.51	0.63	0.37

Table 19. Performance report when the GossipCop dataset was used as “train” and the PolitiFact dataset was used as “test”.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.59	0.56	0.55	0.54	0.55	0.58	0.34
RF	0.59	0.57	0.51	0.41	0.57	0.98	0.02
SVM	0.62	0.61	0.56	0.52	0.60	0.45	0.37

Table 20. Performance report when the PolitiFact dataset was split into “train” and “test”.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.83	0.83	0.81	0.82	0.85	0.14	0.20
RF	0.76	0.76	0.73	0.74	0.84	0.35	0.06
SVM	0.83	0.83	0.81	0.82	0.85	0.17	0.12

Table 21. Performance report when the GossipCop dataset was split into “train” and “test”.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
PA	0.80	0.72	0.72	0.72	0.80	0.15	0.44
RF	0.84	0.82	0.71	0.74	0.84	0.54	0.00
SVM	0.85	0.80	0.75	0.77	0.85	0.35	0.10

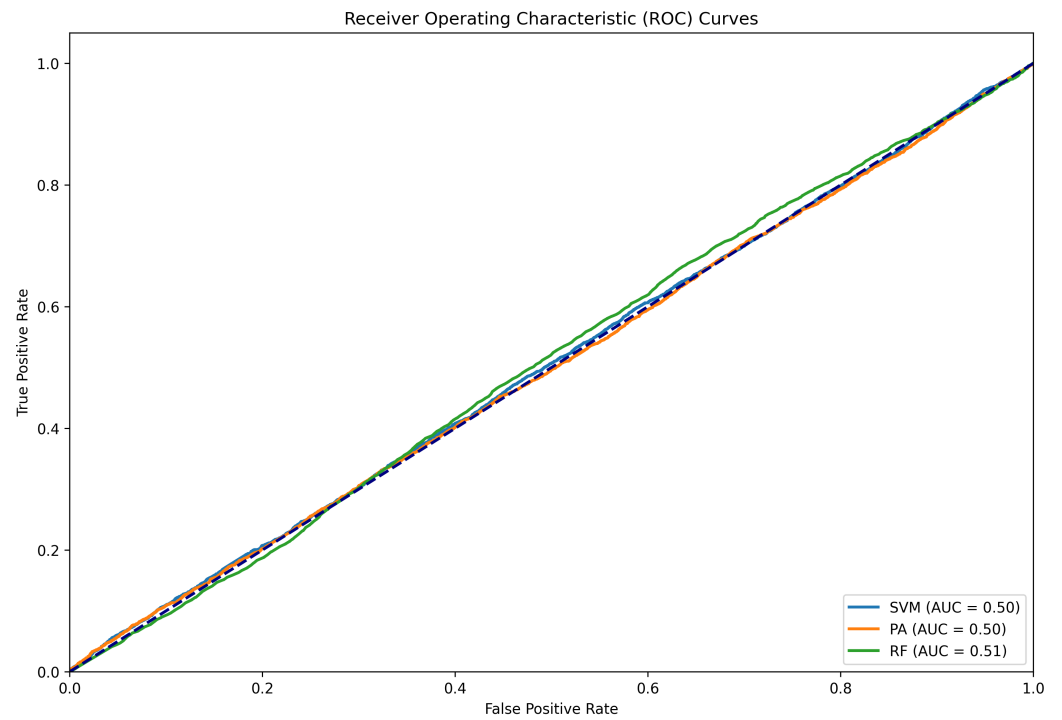


Figure 13. ROC curve when the Politifact dataset was used as “train” and the GossipCop dataset was used as “test”.

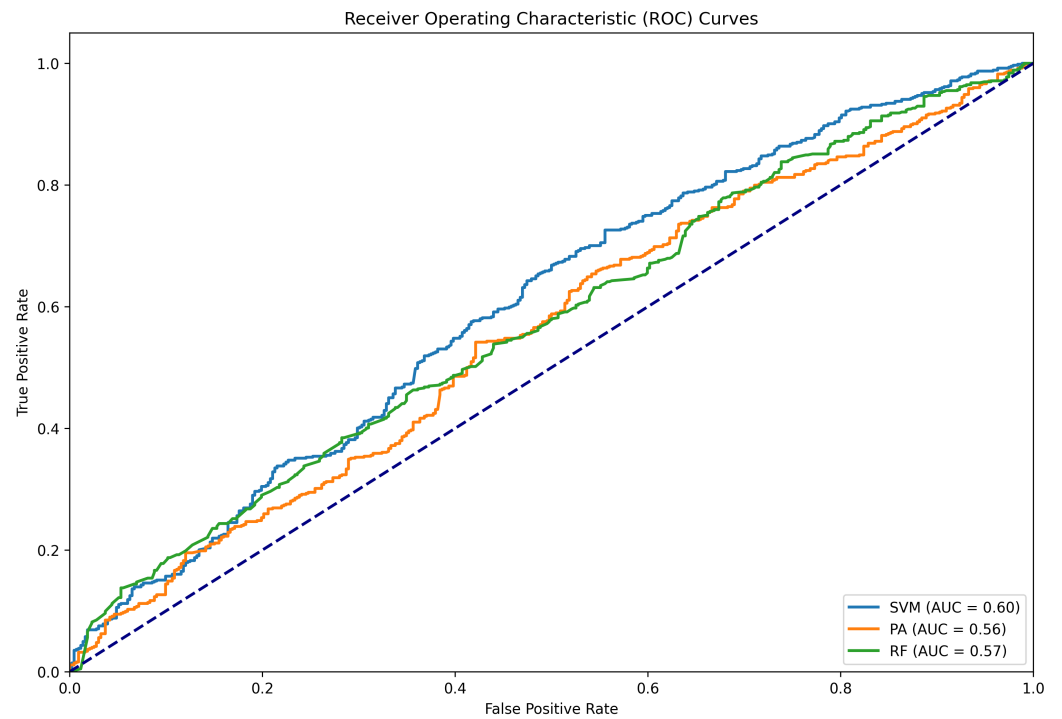


Figure 14. ROC curve when the GossipCop dataset was used as “train” and the Politifact dataset was used as “test”.

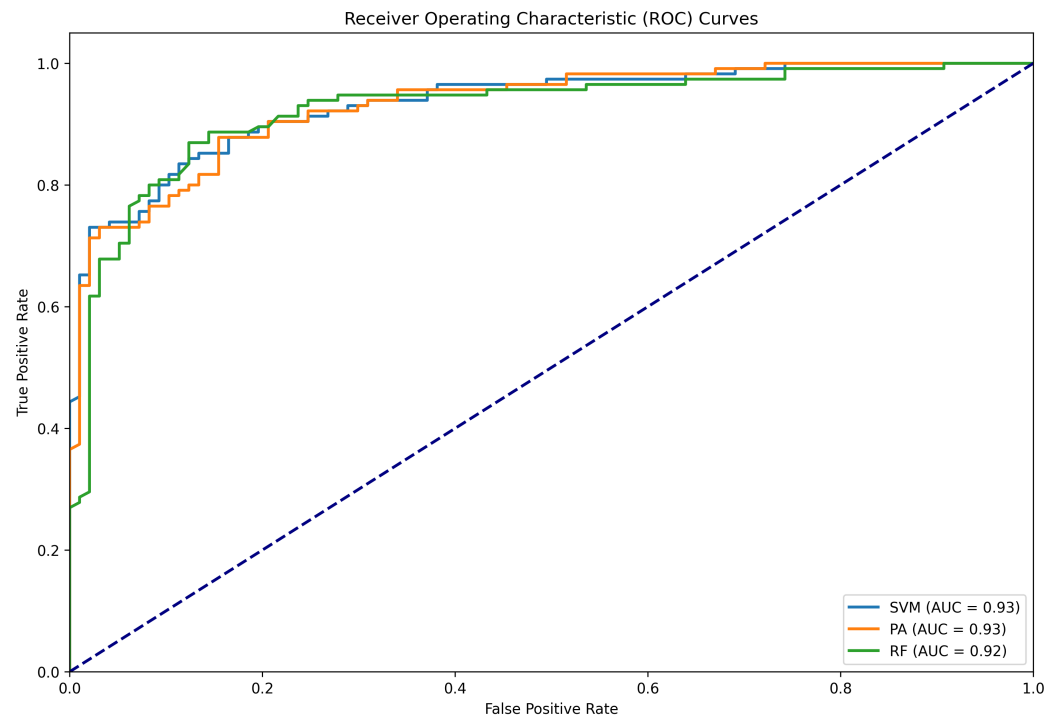


Figure 15. ROC curve when the PolitiFact dataset was split into “train” and “test”.

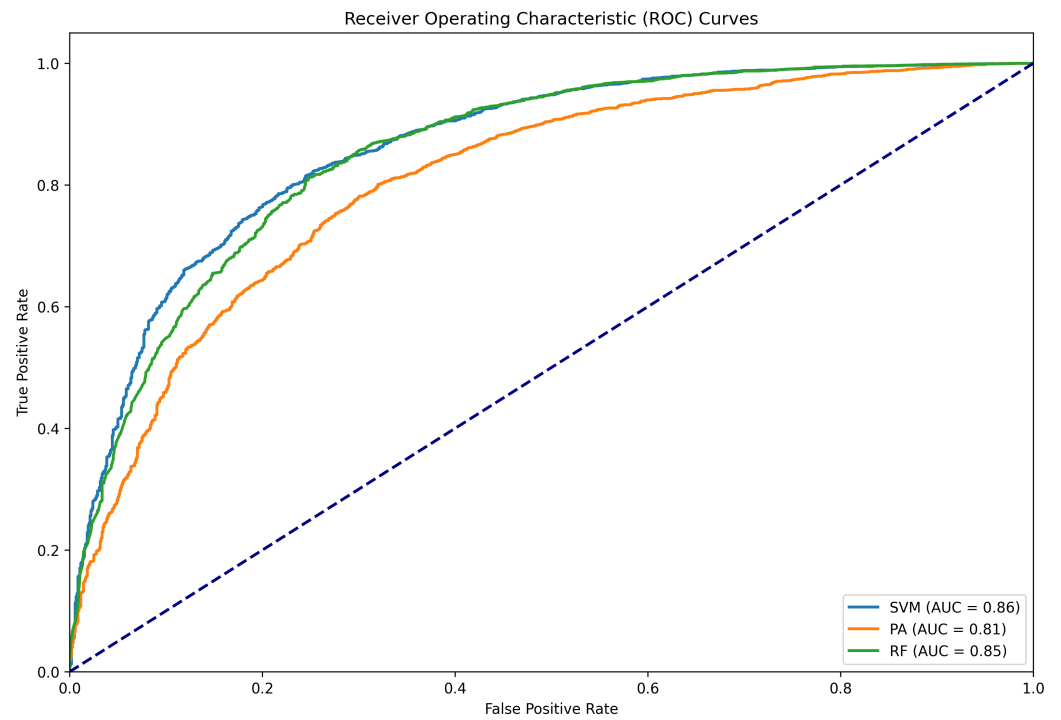


Figure 16. ROC curve when the GossipCop dataset was split into “train” and “test”.

From these tables, we can also observe that the different methods were more or less accurate depending on what was used as “train” or “test” data. From Tables 18 and 19, we can observe that the SVM method performed the best when using the two datasets against each other either as “train” or “test” data. From Table 20, we can observe that the PA and SVM methods achieved identical accuracy scores when splitting the PolitiFact dataset into “train” and “test”. From Table 21, we can observe that the RF and SVM methods were most accurate when splitting the GossipCop dataset into “train” and “test”. However, from these

results, we can observe that the RF method still occasionally struggles with producing good FPR and FNR results.

3.4. Experiment 4: Novel ChatGPT-Generated Dataset

For the fourth and final set of experiments, we used the novel ChatGPT-generated dataset and the LIAR dataset [4]. The novel ChatGPT-generated dataset consists of 200 pieces of false news, generated by ChatGPT-3.5 based on imaginary prompts that we gave to the AI, as well as 100 true news pieces generated by picking actual news articles from various sources, summarizing them, and making ChatGPT-3.5 rewrite the article based on the given summary. As various AI tools have become more widespread in today's world, and seeing how easily and quickly we were able to create false news articles using ChatGPT, it is important to think of ways we could efficiently develop ways to identify AI-generated content. In this experiment, the first classifying experiment was conducted by only using the ChatGPT-generated dataset, the second classifying experiment was conducted using the ChatGPT-generated dataset as the train data and the LIAR dataset as test data, and the third classifying experiment was conducted vice versa.

Like in the previous experiments, we first built a bag-of-words model using CountVectorizer and TfidfTransformer. The methods used for these experiments were the same as in Sections 3.2 and 3.3.

When conducting experiments on the ChatGPT-generated dataset, the split ratio in all of the experiments was 80:20.

The first experiment was conducted only using the ChatGPT-generated dataset. This was implemented to build a baseline of results that we could use to compare the results achieved from the more complex experiments. The results for the first experiment can be seen in Table 22, and the ROC curve for this experiment can be observed in Figure 17. From this table, we can see that all the classifiers were able to achieve good results, but the PA classifier was able to achieve more consistent precision, recall, and F1 scores, as well as low FPR and FNR scores in comparison to the other two classifiers.

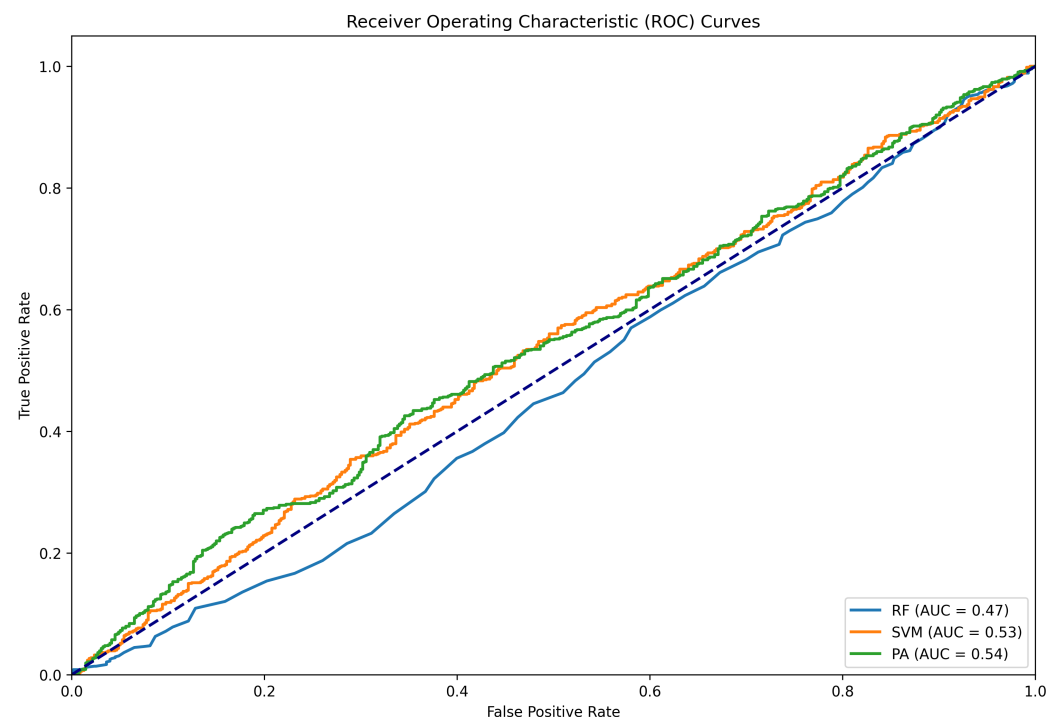


Figure 17. ROC curve for the classifying results conducted using only the ChatGPT-generated dataset.

Table 22. Performance report on the classifying results conducted using only the ChatGPT-generated dataset.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
RF	0.87	0.93	0.56	0.56	0.85	0.00	0.82
SVM	0.84	0.64	0.58	0.60	0.95	0.05	0.32
PA	0.89	0.74	0.79	0.76	0.95	0.05	0.23

The second and third sets of experiments were conducted using the ChatGPT-generated dataset as train data and the LIAR dataset as test data, and vice versa. The results for these experiments can be seen in Tables 23 and 24, and the ROC curves for these experiments can be observed in Figures 18 and 19. The performance of the classifiers dropped significantly, especially when using the ChatGPT-generated dataset as test data. The results obtained from using the ChatGPT-generated dataset as train data were better but lagged behind the results achieved from the first experiment. From these experiments, we can also observe an increase in FPR and FNR scores for all the methods, which stayed relatively low in the first experiment.

Table 23. Performance report on the classifying results conducted using the ChatGPT-generated dataset as train data and the LIAR dataset as test data.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
RF	0.53	0.48	0.49	0.44	0.47	0.75	0.28
SVM	0.56	0.52	0.50	0.43	0.53	0.33	0.63
PA	0.53	0.52	0.52	0.51	0.54	0.26	0.69

Table 24. Performance report on the classifying results conducted using the ChatGPT-generated dataset as test data and the LIAR dataset as train data.

Method	Accuracy	Precision	Recall	F1	AUC	FPR	FNR
RF	0.33	0.42	0.50	0.25	0.47	0.99	0.03
SVM	0.33	0.41	0.47	0.29	0.44	0.78	0.30
PA	0.39	0.47	0.48	0.38	0.49	0.75	0.30

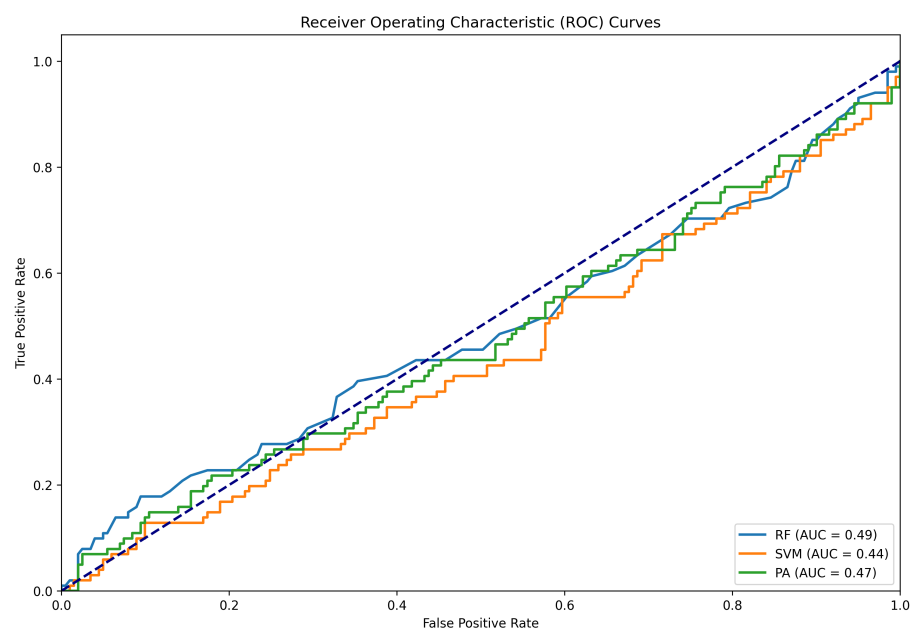


Figure 18. ROC curve for the classifying results conducted using the ChatGPT-generated dataset as train data and the LIAR dataset as test data.

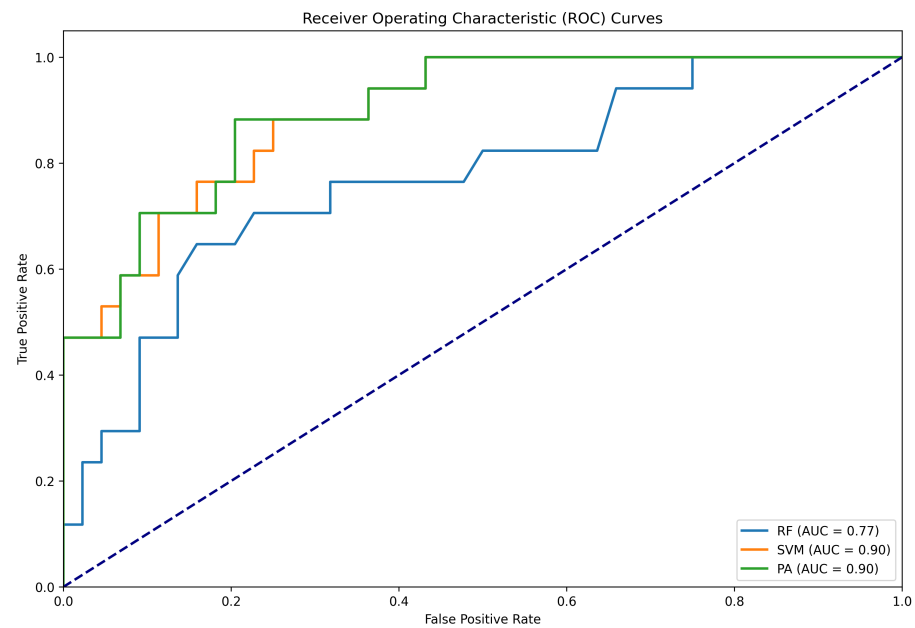


Figure 19. ROC curve for the classifying results conducted using the ChatGPT-generated dataset as test data and the LIAR dataset as train data.

4. Discussion

The results achieved from these experiments are very promising. This paper was able to conclude that linear models and ensemble methods are the most efficient for automatically detecting false news. The three most efficient methods in this research are passive aggressive classifiers, support vector machines, and random forests. These methods not only achieved good accuracy results in many of the experiments but also maintained a stable performance in terms of macro averages of precision, recall, F1-scores, and AUC scores. The relatively good AUC scores indicate that these classifiers have a strong ability to distinguish between false and true news, showing their overall robustness. However, as observed from a lot of the experiments, the RF method struggles to produce good FPR and FNR results.

This paper aimed to compare different types of false news contents and how they differ from each other. In several experiments, classifiers were trained and tested using two different datasets. In all experiments where different datasets were used for training and testing, the accuracy results dropped significantly compared to when the same dataset was used for both training and testing. However, this paper also concluded that the results heavily depend on the specific dataset that was used. For example, in Section 3.3, training using the GossipCop dataset and testing using the PolitiFact dataset were almost 10% higher on all the methods when comparing the results achieved with training using the PolitiFact dataset and testing using the GossipCop dataset. These differences in results could be due to factors such as the size difference between the datasets or the linguistic characteristics of each dataset that make it easier to detect false news in one compared to the other. This is further supported by the observed FPR and FNR scores, which show that the classifiers trained on larger datasets tend to produce fewer false positives and false negatives when tested on smaller datasets, contributing to higher accuracy and AUC scores.

This paper also conducted experiments on whether labeling false news according to different degrees of truthfulness would affect the accuracy of the results. In the Section 3.2, where the LIAR dataset [4] was divided into six different labels based on the degree of falsity or truth, it was found that using multiple labels significantly decreased the accuracy results of the classifiers. When the dataset was re-labeled into just two labels, the classifiers consistently achieved accuracy results of over 50%, the highest accuracy being 91%. In contrast, when the dataset was originally provided with six labels, the classifiers only

achieved accuracy results of about 20%. Additionally, the AUC scores show a similar trend, albeit not as drastic.

Perhaps the most intriguing part of this research was the novel ChatGPT-generated dataset developed specifically for this paper. While building this dataset, we found out how easy it was to generate false news articles in just a matter of a few minutes. We were able to generate about 10 false articles in roughly an hour. The fact that we succeeded in generating false articles that look completely believable easily and quickly shows that there is a need to expand the topic of automated false news detection into the field of detecting computer-generated texts.

In Section 3.4, the focus was on detecting computer-generated false news and understanding how it differs from human-generated false news. The experiments concluded that when using only the ChatGPT-generated dataset, the classifiers were able to achieve similar accuracy results as when using the human-generated datasets. However, when the classifiers were trained using the ChatGPT-generated dataset and tested with the LIAR dataset [4], the achieved accuracy results were around 20% higher than when the training and testing datasets were reversed. The significant difference in accuracy could suggest interesting distinctions between computer-generated and human-generated texts. For instance, ML techniques might be more capable of differentiating false and true from human-generated texts when trained with computer-generated texts than vice versa. This can also be observed from the AUC scores, where the classifiers trained on the ChatGPT-generated datasets achieved higher scores. The FPR and FNR scores also show that the classifiers were less likely to misclassify the human-generated texts when trained on computer-generated data, suggesting nuanced differences in the language used in false news, depending on whether it was created by a computer or a human. These results indicate that detecting computer-generated texts from human-generated ones is a complex task, and they require more sophisticated learning models.

Linguistic Analysis of Language Used in False News

To understand false news better, this study concluded with a simple linguistic analysis of two datasets—the LIAR dataset [4] and the novel ChatGPT-generated dataset. Mainly, we looked at the 20 most common words found in the two datasets, as well as their TF-IDF scores and N-gram frequencies. This way, we can observe the keywords of the datasets, as well as the impact of the words when performing automated detection. From these results, we can also observe whether the language used in false news is different when they are generated by humans or by a computer.

The most common words for the LIAR dataset [4] are presented in a word cloud that can be observed in Figure 20. The most common word that appeared on this dataset was *obama*, with a TF-IDF score of 44.67069 and an N-gram frequency of 271. The TF-IDF scores for other words can be seen in Table 25 and the N-gram frequencies in Table 26.

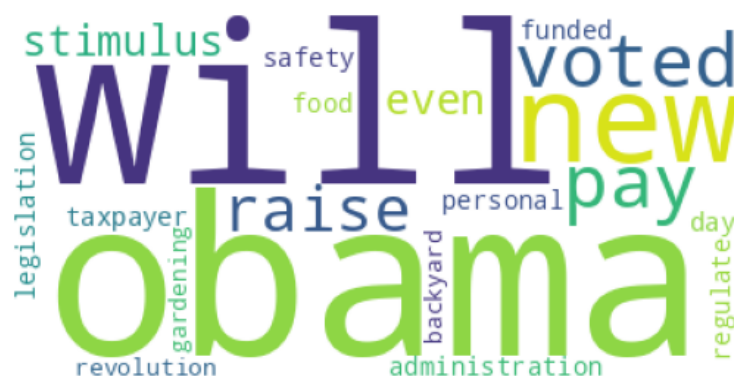


Figure 20. Twenty most common words from LIAR dataset [4].

Table 25. TF-IDF scores of 20 most common words from LIAR dataset [4].

Word	TF-IDF Score
backyard	0.44503
gardening	0.44503
revolution	0.96835
regulate	1.13286
safety	1.66006
personal	2.26182
taxpayer	3.97960
funded	4.05192
food	5.22986
day	5.82952
legislation	6.16053
administration	8.80094
even	9.62488
stimulus	10.08428
raise	10.20364
pay	14.15832
voted	19.63578
new	24.33575
will	26.08082
obama	44.67069

Table 26. N-gram frequencies of 20 most common words from LIAR dataset [4].

Word	N-Gram Frequency
backyard	1
gardening	1
revolution	3
regulate	3
safety	5
personal	8
taxpayer	15
funded	14
food	20
day	26
legislation	27
administration	40
even	45
stimulus	41
raise	40
pay	65
voted	91
new	135
will	158
obama	271

The most common words for the ChatGPT-generated dataset can be observed in Figure 21. The most common word that appeared on this dataset was *potential*, with a TF-IDF score of 3.52128 and an N-gram frequency of 262. The TF-IDF scores for other words can be seen in Table 27 and the N-gram frequencies in Table 28.

When comparing the vocabulary of the two datasets, we can observe that the words appearing in the LIAR dataset [4] are more eye-catching than in the ChatGPT dataset. The words appearing on the ChatGPT dataset tend to be more neutral and do not show as much insight into the actual nature of the generated articles inside the dataset. This observation could argue that using tools like ChatGPT can be very efficient in making false news articles appear more similar in vocabulary to real news articles.



Figure 21. Twenty most common words from ChatGPT-generated dataset.

Table 27. TF-IDF scores of 20 most common words from ChatGPT-generated dataset.

Word	TF-IDF Score
meanwhile	0.27872
serve	0.76586
claim	1.02257
continue	1.28523
challenge	1.15479
remain	1.64038
individual	1.69804
new	1.70962
event	1.71798
future	1.90560
importance	1.98240
may	2.35172
within	2.36911
life	2.41739
public	2.67856
will	3.18160
world	3.38626
human	3.42755
potential	3.52128
incident	3.87482

Table 28. N-gram frequencies of 20 most common words from ChatGPT-generated dataset.

Word	N-Gram Frequency
meanwhile	8
serve	28
claim	35
continue	65
challenge	46
remain	85
individual	69
new	98
event	64
future	113
importance	111
may	154
within	148
life	133
public	167
will	223
world	246
human	183
potential	262
incident	164

5. Conclusions

This study successfully addressed the primary research questions presented in the beginning regarding the effectiveness of various machine learning (ML) methods for automatic false news detection, as well as their performance when compared human-written versus automatically generated news. We discovered the most successful methods for detecting false news by evaluating multiple ML approaches and have gained insights into the subtle differences between news created by humans and news generated by automated systems such as ChatGPT.

Our study not only looks at the strengths and weaknesses of various ML algorithms but also provides an insight into what constitutes false news. We investigated the key characteristics and linguistic patterns that often represent false information, contributing to the larger field of misinformation detection in general.

For the future, we plan to extend our research to include news generated by generative large language models, such as ChatGPT. This future research aims to understand better how news generated by such models differs from human-written news, as well as to assess and build more reliable methods for distinguishing between them. Additionally, we plan to use larger, more comprehensive datasets in addition to more effective and modern language models like BERT [55] or RoBERTa [56].

In conclusion, this study provides a look into what methods are useful for automatically detecting false news and suggests areas for future research in this subject. Our future objectives will refine these techniques, address more modern and complex methods to detect false news, and develop more effective strategies to combat misinformation in the digital age.

The two main contributions of this paper are as follows:

1. Linear and ensemble models work the best in automated false news detection, specifically, passive aggressive classifiers, support vector machines, and random forests.
2. The introduction of a novel ChatGPT-generated dataset for detecting false news created with generative large language models.

Author Contributions: Conceptualization, L.A., M.P. and J.J.; data curation, L.A.; formal analysis, L.A.; investigation, L.A.; methodology, L.A.; project administration, resources, L.A. and M.P.; software, L.A.; supervision, M.P. and J.J.; validation, M.P. and J.J.; visualization, L.A.; writing—original draft preparation, L.A.; writing—review and editing, L.A., M.P. and J.J. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: The original contributions presented in the study are included in the article, further inquiries can be directed to the corresponding author.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Gruener, S. An Empirical Study on False News on Internet-Based False News Stories: Experiences, Problem Awareness, and Responsibilities. Available online: <https://ssrn.com/abstract=3351911> (accessed on 12 August 2024).
2. Hitlin, P. False Reporting on the Internet and the Spread of Rumors: Three Case Studies. *Gnovis Journal*. Georgetown University: Washington, DC, USA, 2003. Available online: <http://pascalfroissart.online.fr/3-cache/2004-hitlin.pdf> (accessed on 12 August 2024).
3. Molina, M.D.; Sundar, S.S.; Le, T.; Lee, D. “Fake news” is not simply false information: A concept explication and taxonomy of online content. *Am. Behav. Sci.* **2021**, *65*, 180–212. [CrossRef]
4. Wang, W.Y. “liar, liar pants on fire”: A new benchmark dataset for fake news detection. *arXiv* **2017**, arXiv:1705.00648.
5. Reis, J.C.S.; Correia, A.; Murai, F.; Veloso, A.; Benevenuto, F. Supervised Learning for Fake News Detection. *IEEE Intell. Syst.* **2019**, *34*, 76–81. [CrossRef]
6. Hsu, T.; Thompson, S.A. Disinformation Researchers Raise Alarms about A.I. *Int. New York Times*, 10 February 2023.
7. Chun, C.-K.; Dong, Y. Misinformation and literacies in the era of generative artificial intelligence: A brief overview and a call for future research. *Emerg. Media* **2024**, *2*, 70–85.
8. Simon, F.M.; Altay, S.; Mercier, H. Misinformation reloaded? Fears about the impact of generative AI on misinformation are overblown: HKS Misinformation Review. *Misinf. Rev.* **2023**. [CrossRef]

9. Apuke, O.D.; Omar, B. Fake news and COVID-19: Modelling the predictors of fake news sharing among social media users. *Telemat. Inform.* **2021**, *56*, 101475. [CrossRef]
10. Waikhom, L.; Goswami, R.S. Fake news detection using machine learning. In Proceedings of the International Conference on Advancements in Computing & Management (ICACM-2019), Jaipur, India, 13–14 April 2019.
11. Thota, A.; Tilak, P.; Ahluwalia, S.; Lohia, N. Fake news detection: A deep learning approach. *SMU Data Sci. Rev.* **2018**, *1*, 10.
12. Gundapu, S.; Mamidi, R. Transformer based automatic COVID-19 fake news detection system. *arXiv* **2021**, arXiv:2101.00180.
13. Svärd, M.; Rumman, P. COMBATING DISINFORMATION: Detecting Fake News with Linguistic Models and Classification Algorithms. 2017. Available online: <https://www.diva-portal.org/smash/get/diva2:1114109/FULLTEXT01.pdf> (accessed on 12 August 2024).
14. Wardle, C. Fake news. It's complicated. *First Draft* **2017**, *16*, 1–11. Available online: <https://firstdraftnews.org/articles/fake-news-complicated/> (accessed on 12 August 2024)
15. Bounegru, L.; Gray, J.; Venturini, T.; Mauri, M. *A Field Guide to "Fake News" and Other Information Disorders: A Collection of Recipes for Those Who Love to Cook with Digital Methods*; Public Data Lab: Amsterdam, The Netherlands, 2018.
16. Chong, M.; Choy, M. An empirically supported taxonomy of misinformation. In *Navigating Fake News, Alternative Facts, and Misinformation in a Post-Truth World*; IGI Global: Hershey, PA, USA, 2020; pp. 117–138.
17. Cambridge-Dictionary. FAKE NEWS | English Meaning. Available online: <https://dictionary.cambridge.org/dictionary/english/fake-news> (accessed on 12 August 2024)
18. Oxford English Dictionary, Oxford-UP. Machine Learning, N. Available online: https://www.oed.com/dictionary/machine-learning_n (accessed on 12 August 2024)
19. Encyclopedia Britannica, Copeland, B. Artificial Intelligence (AI). Available online: <https://www.britannica.com/technology/artificial-intelligence> (accessed on 12 August 2024)
20. Kerner, S.M. What are Large Language models (LLMs)?: Definition from TechTarget. *TechTarget* **2023**, *May 3*, 2–4. Available online: <https://www.techtarget.com/whatis/definition/large-language-model-LLM> (accessed on 12 August 2024)
21. Holdsworth, J.; Scapicchio, M. What Is Deep Learning? *IBM* **2024**. Available online: <https://www.ibm.com/topics/deep-learning> (accessed on 12 August 2024)
22. Oxford English Dictionary, Oxford-UP. Natural Language Processing, N. Available online: https://www.oed.com/dictionary/natural-language-processing_n (accessed on 12 August 2024)
23. Rubin, V.L.; Conroy, N.; Chen, Y.; Cornwell, S. Fake news or truth? Using satirical cues to detect potentially misleading news. In Proceedings of the Second Workshop on Computational Approaches to Deception Detection, San Diego, CA, USA, 17 June 2016; pp. 7–17.
24. Karimi, H.; Roy, P.; Saba-Sadiya, S.; Tang, J. Multi-source multi-class fake news detection. In Proceedings of the 27th International Conference on Computational Linguistics, Santa Fe, NM, USA, 20–26 August 2018; pp. 1546–1557.
25. Oshikawa, R.; Qian, J.; Wang, W.Y. A survey on natural language processing for fake news detection. *arXiv* **2018**, arXiv:1811.00770.
26. Das, A.; Liu, H.; Kovatchev, V.; Lease, M. The state of human-centered NLP technology for fact-checking. *Inf. Process. Manag.* **2023**, *60*, 103219. [CrossRef]
27. Ahmad, I.; Yousaf, M.; Yousaf, S.; Ahmad, M.O. Fake news detection using machine learning ensemble methods. *Complexity* **2020**, *2020*, 8885861. [CrossRef]
28. Wu, Y.; Zhan, P.; Zhang, Y.; Wang, L.; Xu, Z. Multimodal fusion with co-attention networks for fake news detection. In Proceedings of the Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, Online, 1–6 August 2021; pp. 2560–2569.
29. Nadeem, M.I.; Ahmed, K.; Li, D.; Zheng, Z.; Alkahtani, H.K.; Mostafa, S.M.; Mamyrbayev, O.; Abdel Hameed, H. EFND: A Semantic, Visual, and Socially Augmented Deep Framework for Extreme Fake News Detection. *Sustainability* **2023**, *15*, 133. [CrossRef]
30. Singh, V.; Dasgupta, R.; Sonagra, D.; Raman, K.; Ghosh, I. Automated fake news detection using linguistic analysis and machine learning. In Proceedings of the International Conference on Social computing, Behavioral-Cultural Modeling, & Prediction and Behavior Representation in Modeling and Simulation (SBP-BRiMS), Washington, DC, USA, 5–8 July 2017; pp. 1–3.
31. Ahmed, H.; Traore, I.; Saad, S. Detecting opinion spams and fake news using text classification. *Secur. Priv.* **2018**, *1*, e9. [CrossRef]
32. Mitrović, S.; Andreoletti, D.; Ayoub, O. Chatgpt or human? detect and explain. explaining decisions of machine learning model for detecting short chatgpt-generated text. *arXiv* **2023**, arXiv:2301.13852.
33. Tacchini, E.; Ballarin, G.; Della Vedova, M.L.; Moret, S.; De Alfaro, L. Some like it hoax: Automated fake news detection in social networks. *arXiv* **2017**, arXiv:1704.07506.
34. Del Tredici, M.; Fernández, R. Words are the window to the soul: Language-based user representations for fake news detection. *arXiv* **2020**, arXiv:2011.07389.
35. Shu, K.; Mahudeswaran, D.; Wang, S.; Lee, D.; Liu, H. Fakenewsnet: A data repository with news content, social context and dynamic information for studying fake news on social media. *arXiv* **2019**, arXiv:1809.01286. [CrossRef]
36. Liu, X.; Nourbakhsh, A.; Li, Q.; Fang, R.; Shah, S. Real-time rumor debunking on twitter. In Proceedings of the 24th ACM International Conference on Information and Knowledge Management, Melbourne, Australia, 18–23 October 2015; pp. 1867–1870.
37. Breiman, L. Random forests. *Mach. Learn.* **2001**, *45*, 5–32. [CrossRef]

38. Eronen, J.; Ptaszynski, M.; Masui, F.; Smywiński-Pohl, A.; Leliwa, G.; Wroczynski, M. Improving classifier training efficiency for automatic cyberbullying detection with feature density. *Inf. Process. Manag.* **2021**, *58*, 102616. [CrossRef]
39. Lawton, G.; Burns, E.; Rosencrance, L. What Is Logistic Regression?: Definition from TechTarget. TechTarget, 25 April 2024. Available online: <https://www.techtarget.com/searchbusinessanalytics/definition/logistic-regression> (accessed on 12 August 2024)
40. Shivani, N.; Nousheen, S.; Bhavani, P.; Shrivani, P. Fake news detection using logistic regression. *Int. J. Adv. Eng. Manag. IJAEM* **2023**, *5*, 1151–1154.
41. Kanade, V. All You Need to Know about Support Vector Machines. Spiceworks Inc. 2022. Available online: <https://www.spiceworks.com/tech/big-data/articles/what-is-support-vector-machine/> (accessed on 12 August 2024)
42. Yazdi, K.M.; Yazdi, A.M.; Khodayi, S.; Hou, J.; Zhou, W.; Saedy, S. Improving Fake News Detection Using K-means and Support Vector Machine Approaches. *Int. J. Electron. Commun. Eng.* **2020**, *14*, 38–42.
43. Islam, N.; Shaikh, A.; Qaiser, A.; Asiri, Y.; Almakdi, S.; Sulaiman, A.; Moazzam, V.; Babar, S.A. Ternion: An autonomous model for fake news detection. *Appl. Sci.* **2021**, *11*, 9292. [CrossRef]
44. Wikipedia. K-Nearest Neighbors Algorithm—Wikipedia, the Free Encyclopedia. 2023. Available online: <http://en.wikipedia.org/w/index.php?title=K-nearest%20neighbors%20algorithm&oldid=1163707353> (accessed on 19 July 2023).
45. Sidharth. Multi-Layer Perceptron Explained: A Beginner’s Guide—PyCodeMates. 2023. Available online: <https://www.pycodemates.com/2023/01/multi-layer-perceptron-a-complete-overview.html> (accessed on 12 August 2024)
46. Kaur, S.; Kumar, P.; Kumaraguru, P. Automating fake news detection system using multi-level voting model. *Soft Comput.* **2020**, *24*, 9049–9069. [CrossRef]
47. Wikipedia. Decision Tree—Wikipedia, the Free Encyclopedia. 2023. Available online: <http://en.wikipedia.org/w/index.php?title=Decision%20tree&oldid=1165073066> (accessed on 19 July 2023).
48. Patil, D.R. Fake news detection using majority voting technique. *arXiv* **2022**, arXiv:2203.09936.
49. Verma, N. AdaBoost Algorithm Explained in Less Than 5 Minutes—Techynilesh. 2022. Available online: <https://medium.com/@techynilesh/adaboost-algorithm-explained-in-less-than-5-minutes-77cdf9323bfc> (accessed on 19 July 2023).
50. Scikit-Learn. 1.5. Stochastic Gradient Descent. Available online: <https://scikit-learn.org/stable/modules/sgd.html> (accessed on 12 August 2024).
51. Scikit-Learn. 1.1. Linear Models. Available online: https://scikit-learn.org/stable/modules/linear_model.html (accessed on 12 August 2024).
52. Crammer, K.; Dekel, O.; Keshet, J.; Shalev-Shwartz, S.; Singer, Y. Online Passive Aggressive Algorithms. 2006. Available online: <https://www.jmlr.org/papers/volume7/crammer06a/crammer06a.pdf> (accessed on 12 August 2024).
53. Sharma, U.; Saran, S.; Patil, S.M. Fake news detection using machine learning algorithms. *Int. J. Creat. Res. Thoughts IJCRT* **2020**, *8*, 509–518.
54. Ahmed, S.; Hinkelmann, K.; Corradini, F. Development of fake news model using machine learning through natural language processing. *arXiv* **2022**, arXiv:2201.07489.
55. Devlin, J.; Chang, M.W.; Lee, K.; Toutanova, K. 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), Minneapolis, MN, USA, 2–7 June 2019; pp. 4171–4186. [CrossRef]
56. Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *arXiv* **2019**, arXiv:1907.11692.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.