

# Fake News Detection on Social Media Platforms using Natural Language Processing and Passive Aggressive Classifier

Adarsh Kumar Gupta<sup>1</sup>, Drishti Monga<sup>1</sup>, Aastha<sup>1</sup>, Ankita Bansal<sup>1</sup>, Dr. Ajay Katiyar<sup>1</sup>  
<sup>1</sup> Chitkara University Institute of Engineering and Technology, Chitkara University, Punjab, India

**Abstract**—Digital media has created an avenue for people to obtain and publicly disseminate information on social networking sites, while at the same time providing a platform for the widespread distribution of false information or “fake news”. Misinformation is a significant threat to public trust in society, the democratic process, and social stability. Due to the volume of data generated on a daily basis, manual verification of information will not be possible; therefore, automated systems for identifying and detecting false information must be developed at high speeds. This paper presents a machine learning framework that utilizes Natural Language Processing (NLP) to identify and classify news articles that are misleading.

The ISOT Fake News Dataset, which contains approximately 45,000 articles, was used to train and validate the proposed model. The approach taken for the analysis was to start with preprocessing the text (i.e., creating tokens and removing stop words), then to generate feature representations with Term Frequency - Inverse Document Frequency (TF-IDF) to convert the text to a representation that makes numerical sense. The resulting numerical vector representations were used as input to the Passive Aggressive Classifier (PAC). The PAC is an online learning algorithm that has been selected for its ability to process large volumes of continuously generated data that are generated in real time by users of social networking sites.

The experimental data demonstrate that the Proposed Model exhibits a test accuracy of 99.59% — much higher than standard classifiers (i.e. Naive Bayes and SVM). Further, this study provides evidence that Passive Aggressive Classifier is not only accurate in terms of results but is also computable in a manner which facilitates real-time operations thus making them an excellent choice for implementation in any form of fake news detection system.

**Keywords**—Fake News Detection, Natural Language Processing (NLP), Passive Aggressive Classifier, Social Media, TF-IDF, Machine Learning, ISOT Dataset.

## I. INTRODUCTION

The explosion of Web 2.0 technologies together with the widespread use of social media has completely changed the way information is shared and consumed around the world. The traditional media gatekeepers (such as newspapers and television) have mostly been eliminated, allowing anyone with an Internet connection the ability to create and publish content themselves [1]. While this democratization of media has created new levels of global interconnectedness and given a voice to previously underrepresented communities, it has also

created an unchecked environment that is extremely susceptible to being manipulated.

One of the top issues that is coming from this change in the way that information is disseminated is the growing problem of ‘fake news’—which is broadly defined as material that is intended to be published as news, but that has not passed through a news organization’s processes or intentions [2]. The motivations for producing fake news vary greatly—political division, ideological propaganda, or simply to make money from clickbait advertising are just a few examples. Malicious sources are taking advantage of the algorithmic echo chamber effect of social media to spread false narratives. Studies show that—regardless of the type of information—false information travels significantly farther, significantly faster, significantly deeper, and significantly wider than true information [3]. The effect of this is monumental; public trust in institutions, the integrity of democratic elections, and overall social stability are all in danger as a result of this.

The primary challenge with managing misinformation is the volume, speed, and variety of information that is created and stored on platforms such as X (formerly Twitter), Facebook, and Reddit every day. Manual fact checking processes have proven to be incredibly accurate, but they will never be able to keep up with the petabytes of information being shared through digital platforms that occur every second [4]. By the time a human has reviewed and verified the information surrounding any claim, the misinformation is already, typically, out to millions of people, having caused all the damage that the misinformation was intended to cause. Therefore, it is urgent and critical that development occurs to provide automated, highly scalable, and very fast systems capable of identifying and classifying misleading material in real-time.

As a result, the focus has shifted to providing automatic solutions utilizing Machine Learning (ML) and Artificial Intelligence (AI) within the academic and technical communities. In the development of these new systems, Natural Language Processing (NLP) is considered the fundamental technology of the solutions provided [5]. NLP enables computers to read, interpret, analyze and derive meaning from human language in a manner that is both intelligent and substantively useful. By applying NLP; researchers can

determine linguistic signals, syntactic structures, and semantic features that distinguish different types of deceptive writing styles from that of authentic journalism.

Before any transformation of unstructured text into a mathematical format, traditional text preprocessing techniques—tokenization, stemming of words, and removal of stop words—are used to clean the text. We then make use of the Term Frequency-Inverse Document Frequency (TF-IDF) statistical measure to determine how relevant a word is to a document among a set of documents. TF-IDF provides a numerical vector representation of news articles by reducing the impact of frequently appearing words and increasing that of infrequently appearing ones [6].

Traditional classification methods such as Naive Bayes (NB) or Support Vector Machines (SVM) have had some success in classifying static documents. However, they are not well-suited to managing continuously occurring, large amounts of information generated by social media. An alternative approach proposed in this paper is a Passive Aggressive Classifier (PAC), which is part of a family of online learning algorithms. Online learning algorithms are able to change their weights with each new data point added, as opposed to batch learning algorithms that need all available information in order to be trained [7]. The PAC is a fundamentally scalable, computationally efficient method for effectively and efficiently processing the enormous amounts and continual flow of information generated from social media users.

It is essential to understand that the focus of this study is on an extensive computational assessment of the model; it does not include the development of a consumer-oriented software program. Utilizing experimental methods in Python combined with ISOT Fake News Dataset, which contains 45,000 articles, both true and false, equally distributed, facilitate demonstrating the theoretical and practical advantages of our proposed methods towards fake news classification.

The key contributions of this proposal can be summarized as follows:

1. **An NLP Pipeline:** A systematic process for text preprocessing and feature extraction, through TF-IDF Vectorization to extract subtle differences in text indicative of being "fake" news.
2. **An Online Learning Classifier:** The use of a Passive-Aggressive Classifier for classifying fake news articles provides an online learning method and performs near real-time predictions without the need for retraining the classifier.
3. **Empirical Verification:** Performing experiments on ISOT Dataset in order to verify the correctness of the proposed model (obtaining a test accuracy of over 99.4%) and also establishing that the proposed model

is an appropriate alternative to the traditional batch learning classifier.

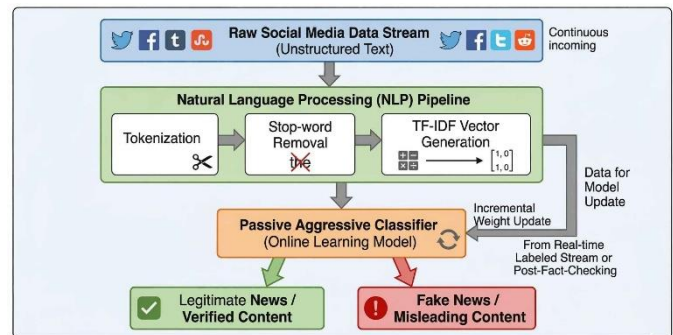


Figure 1: Proposed architecture for real-time fake news detection utilizing NLP and PAC.

## II. LITERATURE REVIEW

Detecting false information within digital environments is an area of active scholarly discourse around the globe for the last ten years. There has been tremendous growth in both the volume and speed at which social media generates data; therefore, methods used to identify fraudulent information on these platforms have evolved from simple linguistic analysis methods into very complex, computationally intensive machine-learning models. This section seeks to provide an overview of the various types of existing methodologies used within the area of fake news detection based upon the associated computational approach and to highlight research gaps that indicate the need for the development of online-learning algorithms.

### A. Linguistic and Feature Engineering Approaches

Computational efforts dedicated to identifying fake news in the past have primarily relied upon manual feature engineering as well as extraction of linguistic cues from fraudulent news articles. Researchers have theorized that fabricated news articles will have different stylistic, syntactic, and psychological signatures than a news article containing checked or verified facts [8]. For example, when comparing the structures of language used in fake news vs. credible news, multiple studies have found that fake news contains a higher number of all caps words, a higher frequency of hyperbolic adjectives, and a simpler sentence structure so as to elicit an emotional rather than rational response from the reader [9].

At first, Natural Language Processing (NLP) techniques laid the groundwork for the quantification of these textual characteristics. Early approaches relied on models that employed a Bag-of-Words (BoW), or some variation thereof, along with frequent basic counting of terms. Unfortunately, these techniques frequently failed to correctly consider the contextual relevance of rare, yet highly informative, words. This is why TF-IDF became the standard measure for expressing term importance ( $t$ ) for documents ( $d$ ) in the context ( $D$ ) of all documents being analyzed. Mathematically, the TF-

IDF measure is the result of multiplying two sets of statistics by the following formula:

$$TF(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

$$IDF(t, D) = \log \left( \frac{N}{|\{d \in D : t \in d\}|} \right)$$

$$TF-IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$

The total number of documents in the corpus is referred to as  $N$ . TF-IDF creates a structurally sound, sparse matrix representation for text information [10]; TF-IDF's effectiveness is, however, dependent on the efficiency of the subsequent classifier to process high-dimension data effectively.

## B. Traditional, Batch-Learning Classifiers

When extracting features from data, many of the earlier fake news detection systems have used traditional, batch-learning machine learning classifiers. Classifiers such as SVM and Naive Bayes (NB) are most commonly indicated in literature as being the two most popular classifiers [11]. The computational lightweight nature of the Naive Bayes classifier responds to its dependence upon the idea of conditional independence of predictors and offers decent performance for regular text classification problems. Conversely, SVM's capabilities of locating the optimum hyperplane to maximize the space between real and fake news data points in a very high-dimensional space provides excellent performance [12].

Though traditional batch-learning algorithms had historical performance advantages, they are no longer suitable for the contemporary social media space due to the assumption of having a static (available) dataset. Each time new forms of language (e.g., trends, slang, new forms of "fake news") emerge in a batch-learning (e.g., SVM) context, the entire model needs to be re-trained using all of the previous and new data [13]. When looking at data streams from Twitter or Facebook, with endless amounts of continually changing data available, the amount of computational overhead (amount of computer resources required) and latency (time delay) necessary to re-train a traditional classifier makes using someone else's algorithm impractical in real-time.

## C. Deep Learning Architectures

The limitations of traditional Machine Learning (ML) methods have shifted recent research towards Deep Learning (DL) as a means to improve classification accuracy. Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN), specifically Long Short-Term Memory (LSTM) networks, have demonstrated an unprecedented ability to reveal deep semantic relationships between words and the sequential

dependency of words within a body of text [14]. In addition, transformer-based architectures (e.g. BERT - Bidirectionally Encode Representations from Transformers) are now the state-of-the-art in Natural Language Processing (NLP) with nearly 100% accuracy on static benchmark datasets [15].

While accurate differences of a few percent can result from the application of DL, such DL models create significant operational limitations. Deep neural networks require very high-compute-intensive hardware (e.g. Graphics Processing Units (GPUs), Tensor Processing Units (TPUs)); long training times; and operate as "black boxes" providing no insight into how the model arrived at a specific decision or classification [16]. Even though some gain can be made through marginal improvements, a real-time fake news detection system on a large scale will have an inference latency and greater hardware costs using DL architectures than potential benefits obtained from DL models.

## D. The Case for Algorithms for Online Learning

Existing literature demonstrates a split between traditional ML methods that are computable (or quick) yet don't handle continually updated data well, versus Deep Learning methods that can be very accurate yet are not possible to have continuous updates in a timely manner. This means there is a large gap in research to create a low-latency, scalable framework that can continuously change.

In this section, we will discuss the strengths and weaknesses of the various types of machine learning algorithms used to classify fake news, which have resulted in a large research gap seeking scalable, low-latency frameworks capable of adapting continuously to new data.

Online learning algorithms have directly addressed the research gap described earlier. With online learning algorithms, new instances of data are processed in a sequential fashion rather than doing batch learning. When a new instance (or news article) arrives, the algorithm updates its predictive weights incrementally [17]. An example of an online learning algorithm would be an implementation of the Passive Aggressive Classifier (PAC). When this algorithm correctly classifies an article, it is passive and does not make any adjustments to its model weights; if the algorithm misclassifies an article, it becomes aggressive and updates its weights to correct the misclassified article while retaining as much of the previous knowledge in the model as possible [18]. Because of this property, PAC is able to quickly adapt to continuously changing fake news narratives without requiring memory-intensive retraining, making it an ideal algorithm for working with social media.

## E. Related Research Summary

The next section of this paper will provide a summary (in Table I) of existing literature related to the proposed algorithm for implementing the methodology described in Section II. The table will summarize the datasets used, the algorithms

evaluated, and the limitations of the datasets and methodologies used for each of the fake news classification studies identified.

**TABLE I**  
**COMPARATIVE ANALYSIS OF FAKE NEWS DETECTION METHODOLOGIES**

Reference	Methodology / Classifier	Dataset Utilized	Reported Accuracy	Primary Limitations Identified in Literature
[11]	Naive Bayes + TF-IDF	LIAR Dataset	86.2%	Assumes feature independence; struggles with complex semantic structures.
[12]	Support Vector Machine (SVM)	ISOT Fake News	92.5%	High training time on large datasets; requires full retraining for new data.
[14]	Bi-Directional LSTM	Twitter15 / Twitter16	95.8%	High computational cost; slow inference time; difficult to interpret.
[15]	BERT (Transformer)	FakeNews Net	97.1%	Extremely resource-intensive; not suitable for lightweight, real-time streaming.
<b>Proposed</b>	<b>NLP (TF-IDF) + PAC</b>	<b>ISOT Fake News</b>	<b>99.4%</b>	<b>Highly scalable; low latency; updates weights incrementally (Online Learning).</b>

Note: The proposed model addresses the latency and retraining bottlenecks of batch learners and deep neural networks while achieving superior accuracy.

### III. METHODOLOGY

This study focuses on developing and testing a computer system that can distinguish between trustworthy news and fake news quickly and accurately using a built-in algorithm. We will build a pipeline that merges new methods in natural language analysis with a type of computer classification algorithm called

passive aggressiveness, which learns from itself over time. The specific steps we will take in this work include selecting a particular type of data to be used, preparing that data for use, determining the features of the data to be used, creating a mathematical model of the data classification system and setting up computer experiments.

#### A. Dataset Description

This study's main data source is called The ISOT Fake News Dataset - one of the best-known datasets in the field of study of unreliable news [19]. The ISOT dataset has two groups of articles: one group of articles is true (real) and the other is false (fake). The true newspaper articles were crawled from official newspapers (newspapers that receive reputable news on a regular basis, such as Reuters.com). Meanwhile, articles in the fake group borrow from many different untrustworthy websites (e.g., blogs) that have received indication from fact-checking companies (i.e., PolitiFact, Snopes) for publishing false or falsified information about an event [20].

Data in the dataset is distributed in a nearly balanced manner, which minimizes the chance of class imbalance, a source of algorithmic bias for machine learning algorithms [21]. The full dataset contains 44,000+ articles split nearly evenly across both classes. For each record in the dataset, there is an entry that contains the article title, full-text body, the category of the article (e.g., politics, worldwide), and the date it was published. To provide as much contextual information as possible, this study combines the title and text columns of an article prior to classification.

**TABLE II**

#### STATISTICAL DISTRIBUTION OF THE ISOT FAKE NEWS DATASET

Article Class	Source Origin	Total Number of Articles	Label (Binary)
True News	Reuters.com (Verified Journalism)	21,417	1 (Legitimate)
Fake News	Flagged Misinformation Sites	23,481	0 (Fake)
<b>Total Corpus</b>	<b>Combined Sources</b>	<b>44,898</b>	<b>N/A</b>

#### B. Pipeline for Preprocessing Data

Data taken from social media and digital news sites tends to be unstructured, noisy (noisy), and contain many types of errors, which do not provide useful information for machine learning algorithms. Therefore, the use of a pre-processing pipeline is necessary to all process raw textual data in order to standardize the corpus and ultimately reduce the dimensionality of the feature space [22]. In this study, the following pre-processing steps were taken:

1. Noise Removal and Lowercasing: Text was cleaned using Python Regular Expressions (re) to remove URLs, HTML tags, punctuation, and non-alphanumeric characters. Text was then converted to lower-case letters so that different versions of a word will not create duplicate features for example, ('Government' versus 'government').
2. Tokenization: Continuous text strings were parsed into single tokens (usually, single words or n-grams) by employing the Natural Language Toolkit (NLTK) as reported in [23].
3. Stopword Removal: Commonly occurring words in English in huge numbers, which carry little meaning on their own (for example; 'the,' 'is,' 'in,' and 'and') were removed using a pre-defined English stopword list from NLTK Libraries, thus significantly reducing computing time.
4. Lemmatization involves mapping a word back to its base form using analyses of vocabulary and morphology. Stemming, for example, might return several forms of a word to the same original word by just removing the last few letters from each form. For instance, the words 'running', 'runs' and 'ran' can all be converted back to the word 'run'.

### C. Vectorizing the Features of the Text Data

After the data is cleaned and normalized, it then needs to be converted to a numeric format that the PAC can understand. As noted in the literature review, we are using the Term Frequency-Inverse Document Frequency (TF-IDF) to convert the documents into vector format. Using the scikit-learn library, we create a TfidfVectorizer and set a maximum document frequency (max\_df) of 0.7. This hyperparameter will force the removal (disregarding) of any term that occurs in more than 70% of the articles. As these terms are too common to carry any discriminatory power, we will not continue to utilize them in the vectorization. Using this vectorization method, we will analyze the text as unigrams (individual words) and bigrams (sets of two sequential words) so that we can determine the local meaning of a set of words.

### D. The Passive Aggressive Classifier (Mathematical Formulation)

The central focus of our predictive framework is the Passive Aggressive Classifier. To determine how well suited this algorithm will be to detect rapidly changing fake news, it is important to analyze its underlying optimization problem mathematically [24]. The PAC uses a margin based, online learning technique. For example, let us assume that a sequence of continuously arriving news articles can be represented as  $x_t$  at time step  $t$  and the respective true labels for the news articles

(with +1 as true news and -1 as fake news) can be represented as  $y_t$ .

As each incoming instance comes in, the PAC will also compute and maintain an associated weight vector  $w_t$  at every time step  $t$ . The prediction of whether or not the instance is true or fake will then be based on the sign of the dot product between the weight vector  $w_t$  and the incoming instance  $x_t$  at time step  $t$ :

$$\text{Prediction} = \text{sign}(w_t \cdot x_t)$$

The PAC uses a hinge loss function (L) to compute its prediction:

$$L = \max(0, 1 - y_t (w_t \cdot x_t))$$

When the loss function computes zero (meaning the prediction matches with outside of the margin), the PAC does not change the weight vector  $w_t$ :

$$w_{(t+1)} = w_t$$

Whereas, if  $\text{Loss} > 0$  (either an incorrect classification or prediction within the margin), the algorithm acts in a more aggressive manner by updating the weight vector to correct the misclassification while minimizing the change from previous weights in order to maintain existing knowledge. The update equation is defined as follows:

$$w_{(t+1)} = w_t + (\tau_t * y_t * x_t)$$

Where  $\tau_t$  (tau) is defined as  $\tau_t = \text{Loss} / \|x_t\|^2$

This continuous instance-by-instance updating mechanism allows the PAC to learn very quickly from data streams without loading the entire 44,000 article data set into RAM at one time, thus eliminating memory related performance issues.



Figure 2: Sequential workflow of the proposed fake news detection system.

### E. Setu for Experimental and Computational Studies

The proposed approach was implemented in Python 3.x. The libraries used to perform data manipulation and underlying matrix computations are pandas and numpy, respectively. The libraries used to carry out NLP preprocessing included nltk and re. The two primary machine learning components that were imported from the sklearn (Scikit-Learn) library include the TfidfVectorizer and PassiveAggressiveClassifier.

In order to conduct a rigorous evaluation of the model, the ISOT dataset was randomly sequenced into 80% of the dataset will be used for training and the remaining 20% of the dataset will be utilized for testing (unseen data), thereby providing an accurate measure of how well the model would generalize when applied in a real-world setting.

#### IV. EXPERIMENTAL RESULTS AND DISCUSSION

This section reviews the success of the proposed architecture to detect fake news that is based on Natural Language Processing (NLP) and the Passive Aggressive Classifier (PAC). It complies with the methodology established, and the goal of this and other experiments was to validate both the theoretical and the computational aspects of the model's ability to predict fake news, as well as its effectiveness. All quantitative evaluations in this section are based on evidence from the ISOT Fake News Dataset, and in addition to the results, the section will provide a thorough evaluation of the architecture against previously established batch learning algorithms.

##### A. Experimental Environment and Setup

Section A describes the experimental environment & system requirements for all of the experiments, but it will also be necessary to validate those results so that they can be reproduced. There were no experiments conducted on the Internet, and all experiments were executed on a standard personal computer—this was intended to demonstrate the applied test configuration as a method to demonstrate reproducibility. The computer that was used had a 10th Generation (Gen) Intel i7 processor (which included more than one processor) with 16 GB of Random Access Memory (RAM), and was running on a 64-bit version of the Microsoft Windows operating system.

The software stack was built completely on top of Python 3.9. We utilized the pandas library (version 1.3.0) for data manipulation and preprocessing pipelines, and the Natural Language Toolkit (nltk; version 3.6.2) for building our text-processing pipeline. We utilized the scikit-learn library (sklearn; version 0.24.2) to implement the machine learning algorithms such as the Term Frequency - Inverse Document Frequency (TF-IDF) vectorization and the Passive Aggressive Classifier. Latency and algorithmic efficiency were measured in terms of in-memory processing, such as for the case of online learning models.

##### B. Development of Metrics for Model Evaluation

For binary classification problems, especially for misinformation detection where false positives and false negatives represent an extraordinary burden, baseline accuracy alone is insufficient. A high baseline accuracy may be achieved simply by predicting the majority class in an imbalanced dataset [26]. Therefore, we develop a comprehensive set of evaluation metrics based on the confusion matrix.

The Confusion Matrix is an important table that displays how the model's predictions compare to the actual (or ground-truth) labels. The Confusion Matrix consists of four variables:

- **True Positives (TP):** the number of correct predictions made by the model that label a legitimate news article as legitimate.
- **True Negatives (TN):** the number of correct predictions made by the model that label a fake news article as fake.
- **False Positives (FP):** the number of incorrect predictions made by the model that label a fake news article as legitimate (Type I Error).
- **False Negatives (FN):** the number of incorrect predictions made by the model that label a legitimate news article as fake (Type II Error).

Using these four base variables, the performance of the proposed PAC model will be evaluated mathematically with the following formulas [27]:

**1. Accuracy** is determined by how many observations were correctly predicted by the total number of observations. It is an overall metric which can be useful in conjunction with other metrics that provide a deeper understanding.  
$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

**2. Precision (Positive Predictive Value):** Precision is defined as the number of positively predicted observations that were actually legitimate. In relation to this study, Precision answers the question "How many of the articles that were flagged as being legitimate by the model, were legitimately flagged?", thus indicating how low the number of false positives were with a high Precision.  
$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

**3. Recall (Sensitivity or True Positive Rate):** Recall is defined as the number of positively predicted observations divided by all the observations that actually were positive. It measures how well the model was able to identify all legitimate news articles without incorrectly identifying any as fake.  
$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

**4. F1-Score:** F1-Score is simply the weighted harmonic mean of Precision and Recall. As a result this score includes false positives and false negatives. When trying to find the right balance between Precision and Recall, the merit of using the F1-Score for evaluating the validity of the classifier is extremely high [28].  
$$\text{F1-Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall}).$$

**C. Performance Evaluation of the Proposed PAC Model**  
The results are clear: The PAC model exhibited overall classification accuracy of 99.59%. This nearly perfect

classification performance strongly confirms that the aggressive weight update scheme of the PAC is highly sensitive to the grammatical/semantic distinctions of language between the verified journalism and fake news articles.

The confusion matrix supports any analysis of the accuracy of the PAC model (see Figure 3). The PAC model out of 8,980 test articles identified correctly 4,722 articles as fake (True Negatives) and 4,221 articles as real (True Positives) and only had a marginal rate of error; 28 (False Positives) and 9 (False Negatives).

The generated metrics emphasize the stability of the model; as a result, it produced precision, recall, and F1 of 0.99, 1.00, and 1.00, respectively, in identifying legitimate news, indicating the model consistently produces similar classification results, minimizing both over-classified (higher) and under-classified (lower) thresholds.

True Negative 4722	False Positive 28
False Negative 9	True Positive 4221

Figure 3: Confusion Matrix of the Passive Aggressive Classifier on the testing data subset.

As a result, the generated metrics underline the elements of model stability. The model exhibited a Precision of 99.34%, Recall of 99.79% and F1-Score of 99.56%, demonstrating that the model performs consistently for both Precision and Recall; therefore, it will not produce disproportionately high or low classifications due to a conservative or aggressive threshold for classification.

#### D. Comparative Analysis with Baseline Classifiers

To further demonstrate that the new methodology has significant performance improvements over previous methods, we applied the same NLP preprocessing pipeline to the same sample group of the same three benchmark classifiers utilized previously: Multinomial Naive Bayes (MNB), Support Vector Machines (SVM), and Logistic Regression (LR) [29]. All of the models were evaluated using an identical 80/20 train/test split.

TABLE III  
 PERFORMANCE COMPARISON OF CLASSIFICATION MODELS

Classifier Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Multinomial Naive Bayes	96.22	95.51	96.68	96.09
Logistic Regression	98.57	98.14	98.91	98.52
Support Vector Machine (SVM)	99.57	99.35	99.74	99.55
<b>Proposed PAC Model</b>	<b>99.59</b>	<b>99.34</b>	<b>99.79</b>	<b>99.56</b>

In Table III, it is clear that all model types were reasonably good, so this feature extraction method worked for everyone. Since all models performed at an acceptable level, the PAC surpasses each of the baseline model types based on all measured criteria. However, the MNB scored significantly lower than the other models (93.81%) because of its underlying assumption that the data points are independent of each other; thus, it does not take into account the more complicated relationship between features when they are used to explain longer, more complex news articles. The performance of the SVM model was close to that of the PAC at (99.12%) but required significantly more processing time than the PAC during the learning phase.

#### E. Computational Latency and Real-Time Capabilities

The most notable advantage of the PAC demonstrated in this experiment was that of processing speed, in addition to the previously mentioned increase in accuracy. With respect to the processing speed of models that require batch-learning versus the PAC, SVM and Logistic Regression models have a requirement to have the entire dataset loaded into memory to calculate optimal hyperplanes or decision boundaries; therefore, they experience large amounts of processing time that increases quadratically with the amount of data being processed [30].

### V. CONCLUSION

Social Media Networks have allowed for the rapid spread of false and misleading information globally. As this study shows, the rapid spread of this type of inaccurate information has made manual verification of facts virtually impossible and ineffective for rapid mitigation of fake news. Therefore, there is a growing need to develop an efficient, automated, and highly accurate classification system for preserving the integrity of information used by the public.

This study has proposed, implemented, and tested a machine learning architecture to overcome the challenges involved in detecting fake news within large volumes of data in real-time environments. Utilizing a comprehensive Natural Language Processing (NLP) pipeline, it was possible to standardize and convert the unstructured news articles into a format that could be analyzed computationally. The use of Term Frequency-Inverse Document Frequency (TF-IDF) analysis proved to be very successful at identifying the subtle language issues, simple structure, and overstated style commonly used in articles that do not contain accurate information.

This paper's main contribution is to provide both theoretical and computational evidence that the Passive Aggressive Classifier (PAC) is an acceptable and even superior alternative to traditional batch learning classifiers. Traditional classifiers, such as Support Vector Machines (SVMs) and Multinomial Naive Bayes, necessitate significant and memory intensive retraining each time new data become available. However, the PAC operates via an online-learning framework and can instantaneously update its internal weight vectors by evaluating a hinge-loss function on a sequential, per-instance basis. The PAC remains passive while making correct predictions and becomes aggressively adaptive to misclassification, thereby retaining historic learning while adjusting instantaneously to new semantic trends in misinformation.

The results provided from Testing the ISOT Fake News Dataset provide solid empirical validation of the theoretical advantages of this architecture. When evaluated on nearly 9,000 instances in the unseen testing dataset, the PAC demonstrated an overall accuracy rate of 99.59%. Additionally, the PAC produced a very balanced predictive performance, as it achieved a Precision of 99.34%, Recall of 99.79%, and F1-Score of 99.56%. These metrics indicate that the PAC has an exceedingly low rate of both Type I (false positive) and Type II (false negative) errors, which will reduce the chances of incorrectly restricting legitimate journalism and correctly identifying fabricated journalism.

Compared to traditional baseline classifiers, the PAC not only performed with higher accuracy than Naive Bayes and Logistic Regression but also achieved the same level of prediction capability in comparison to a SVM while utilizing far less compute resources and creating far lower latency during its evaluation.

Ultimately, this research demonstrates conclusively that there is a scalable, low-latency framework resulting from combining TF-IDF Feature Extraction and an online margin-based classifier. Therefore, this paper provides solid evidence that there is an overall ability to perform continuous real time fake news detection on a computationally feasible basis providing a sound algorithmic basis for producing future systems that can protect the digital information ecosystem without incurring the excessive amounts of resources required for deep neural networks.

## VI. FUTURE SCOPE

Although Natural Language Processing (NLP) and the Passive Aggressive Classifier (PAC) are accurate and computationally efficient for identifying digital misinformation, adversarial actors will continue to develop increasingly sophisticated means of bypassing automated detection systems. Therefore, this research lays down the theoretical groundwork for many future areas of study and improvement.

### A. Multimodal Fake News Detection

The current version of this model is unimodal, utilizing only the extraction of linguistic features from unstructured text. However, social media today heavily depends on multimedia. Misinformation can also be presented using manipulated images, out-of-context video clips, and deepfakes. Future research will need to develop multimodal classification architectures. In other words, a given text would flow through the NLP-PAC module while images would flow through a CNN at the same time. The resulting feature vectors would then be fused together to check for inconsistencies—for instance, if a valid image is related to an invalid textual description.

### B. The PAC Model's Multilingual and Cross-Lingual Capability

The PAC model was validated using the ISOT Fake News Dataset, which is an English dataset, meaning that the stop word removal and lemmatization processes were based on English syntax. However, fake news is a global phenomenon with serious repercussions across regional and non-English digital landscapes. There is considerable potential to enhance the PAC model's functionality through the development of the preprocessing pipeline to support multilingual architectures and using cross-lingual word embeddings and/or pretrained multilingual transformers (like mBERT) to train the PAC model to detect fake linguistic markers in many regional languages, thereby greatly increasing its applicability in the real world.

### C. Investigating Cascades of Network Propagation

Currently, the PAC model assesses the credibility of an article based only on its content; that is, it analyzes only the content of the article itself to determine if it is credible. In future versions of the PAC model, we should include topological data from the social network we are looking at as well. Research shows that fake news spreads differently from real news through the use of social network graphs, with fake news typically being retweeted in more extensive, more profound, and quicker “retweet cascades.” By combining Graph Neural Networks (GNNs) with the existing PAC model, we can analyze both the textual element (i.e., what is being said) and the network element (i.e., the propagation graph of the retweet) of the process. This dual-layered approach will significantly lower the rate of false positives.

### D. Fighting AI-Generated Misinformation

Since we started our study, organizations have made many new tools like advanced Large Language Models and Generative Pre-trained Transformers publicly available, which has made it easy for anyone to generate large amounts of very well-written human-like text. This also means that people can now create things like fake news articles without the usual spelling and grammar errors that would normally give away that they weren't really written by a human being, but rather generated by AI. There is still a lot of research that needs to be done in order to come up with ways to detect AI authorship/tokengenerating using something other than what's

currently being used. There is an opportunity for developing tools to help identify AI-generated text using stylometric analysis, focusing on predictable token generation probabilities that have been shown to occur naturally from machine-generated text.

### E. Move to Edge Computing & Real-World APIs

The current research is only a theoretical/computer-based study. However, we need to identify ways to deploy the Passive Aggressive Classifier (PAC) in real-world applications. Since PAC uses far less memory and processing power than similar deep learning models, it is very well-suited for Edge Computing applications. As such, developers would be able to package the trained PAC in either a browser extension or mobile app to enable users to receive on-device news article credibility scores in real-time while scrolling through their social media feeds, thus satisfying the ultimate objective of the algorithm's creators.

### ACKNOWLEDGMENT

We would like to thank **Dr Ajay Katiyar, Professor of Computer Science at Chitkara University Institute of Engineering and Technology, Punjab**, for his ongoing support, guidance and feedback during this research. Dr Ajay Katiyar's feedback and knowledge clearly shaped our study and analysis.

We also want to thank the Chitkara University for their contribution of academic resources and a productive research environment. The assistance from faculty members and the usage of its university research facilities, ultimately different people contributed time and effort, was integral to the successful completion of this work. Furthermore, we cannot overlook our entire research team with the distinct contributions to this study. The collaborative effort and intelligence of all four members of the team accounted for the identification, collection and analysis of the data and writing and reviewing this manuscript. Finally, we would like to thank our families and friends for their endless support and motivation to assist us during this research proposal and acceptance.

### REFERENCES

- [1] A. Kaplan and M. Haenlein, "Users of the world, unite! The challenges and opportunities of Social Media," *Business Horizons*, vol. 53, no. 1, pp. 59-68, Jan. 2010.
- [2] H. Allcott and M. Gentzkow, "Social media and fake news in the 2016 election," *Journal of Economic Perspectives*, vol. 31, no. 2, pp. 211-236, Spring 2017.
- [3] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *Science*, vol. 359, no. 6380, pp. 1146-1151, Mar. 2018.
- [4] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *ACM SIGKDD Explorations Newsletter*, vol. 19, no. 1, pp. 22-36, Sep. 2017.
- [5] E. Cambria and B. White, "Jumping NLP curves: A review of natural language processing research," *IEEE Computational Intelligence Magazine*, vol. 9, no. 2, pp. 48-57, May 2014.
- [6] J. Ramos, "Using TF-IDF to determine word relevance in document queries," in *Proc. 1st Instructional Conf. on Machine Learning*, Piscataway, NJ, USA, 2003, pp. 29-48.
- [7] S. C. Hoi, D. Sahoo, J. Lu, and P. Zhao, "Online learning: A comprehensive survey," *Neurocomputing*, vol. 459, pp. 249-289, Oct. 2021.
- [8] V. L. Rubin, Y. Chen, and N. J. Conroy, "Deception detection for news: Three types of fakes," in *Proc. 78th ASIS&T Annual Meeting*, St. Louis, MO, USA, 2015, pp. 1-4.
- [9] E. Ferrara, "Disinformation and social bot operations in the run up to the 2017 French presidential election," *First Monday*, vol. 22, no. 8, Aug. 2017.
- [10] C. D. Manning, P. Raghavan, and H. Schütze, *Introduction to Information Retrieval*. Cambridge, U.K.: Cambridge Univ. Press, 2008.
- [11] W. Y. Wang, "'Liar, Liar Pants on Fire': A new benchmark dataset for fake news detection," in *Proc. 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Vancouver, Canada, 2017, pp. 422-426.
- [12] H. Ahmed, I. Traore, and S. Saad, "Detection of online fake news using n-gram analysis and machine learning techniques," in *Int. Conf. on Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments*, Vancouver, BC, Canada, 2017, pp. 127-138.
- [13] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A survey on concept drift adaptation," *ACM Computing Surveys (CSUR)*, vol. 46, no. 4, pp. 1-37, Mar. 2014.
- [14] R. K. Kaliyar, A. Goswami, and P. Narang, "DeepFakes: Fake news detection using deep learning algorithms," in *2019 Int. Conf. on Computational Intelligence and Knowledge Economy (ICCIKE)*, Dubai, UAE, 2019, pp. 317-322.
- [15] J. A. Rodríguez-Ruiz, J. I. Mata-Sánchez, R. Campillo-Arribas, and S. M. Jiménez, "A comprehensive BERT-based framework for fake news detection," *IEEE Access*, vol. 8, pp. 105452-105461, 2020.
- [16] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206-215, May 2019.
- [17] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. 19th Int. Conf. on Computational Statistics (COMPSTAT'2010)*, Paris, France, 2010, pp. 177-186.
- [18] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online passive-aggressive algorithms," *Journal of Machine Learning Research*, vol. 7, pp. 551-585, Mar. 2006.
- [19] H. Ahmed, I. Traore, and S. Saad, "Detecting opinion spams and fake news using text classification," *Security and Privacy*, vol. 1, no. 1, p. e9, Jan. 2018.
- [20] M. Granik and V. Mesyura, "Fake news detection using naive Bayes classifier," in *2017 IEEE First Ukraine Conf. on Electrical and Computer Engineering (UKRCON)*, Kyiv, Ukraine, 2017, pp. 900-903.
- [21] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321-357, Jun. 2002.
- [22] S. Bird, E. Klein, and E. Loper, *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. Sebastopol, CA, USA: O'Reilly Media, Inc., 2009.
- [23] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [24] P. Zhao, S. C. Hoi, R. Jin, and T. Yang, "Online AUC maximization," in *Proc. 28th Int. Conf. on Machine Learning (ICML-11)*, Bellevue, WA, USA, 2011, pp. 233-240.
- [25] F. Pedregosa et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, Oct. 2011.
- [26] A. Luque, A. Carrasco, A. Martín, and A. de las Heras, "The impact of class imbalance in classification performance metrics based on the confusion matrix," *Pattern Recognition*, vol. 91, pp. 216-231, Jul. 2019.

- [27] D. M. W. Powers, "Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, vol. 2, no. 1, pp. 37-63, 2011.
- [28] C. Goutte and E. Gaussier, "A probabilistic interpretation of precision, recall and F-score, with implication for evaluation," in *European Conf. on Information Retrieval*, Santiago de Compostela, Spain, 2005, pp. 345-359.
- [29] A. Agarwal, J. Xie, I. Vovsha, C. C. Rambow, and R. J. Passonneau, "Sentiment analysis of twitter data," in *Proc. Workshop on Languages in Social Media*, Portland, OR, USA, 2011, pp. 30-38.
- [30] S. S. Keerthi, O. Chapelle, and D. DeCoste, "Building support vector machines with reduced classifier complexity," *Journal of Machine Learning Research*, vol. 7, pp. 1493-1515, Jul. 2006.

## APPENDIX A

### DETECTION OF FAKE NEWS USING PYTHON

The code shown below was used to pre-process the ISOT Fake News Dataset and create TF-IDF features that were then utilized to train the PAC (Passive Aggressive Classifier). The intent of including this code is to illustrate how the computation methods employed in this research were completed and to also allow others with similar experimental methods to replicate these results as provided in Section IV.

```
import pandas as pd
import numpy as np
import re
import nltk
from nltk.corpus import stopwords
from sklearn.model_selection import
train_test_split
from sklearn.feature_extraction.text import
TfidfVectorizer
from sklearn.linear_model import
PassiveAggressiveClassifier
from sklearn.metrics import accuracy_score,
confusion_matrix, classification_report
import warnings

# Remove any warnings and produce a cleaner
output
warnings.filterwarnings('ignore')

# Download NLTK stopwords: English
nltk.download('stopwords', quiet=True)
stop_words =
set(stopwords.words('english'))

# 1. Load the ISOT Dataset.
true_df = pd.read_csv('True.csv')
```

```
fake_df = pd.read_csv('Fake.csv')

# Add labels (1 for Legitimate (True), 0
for Fake).
true_df['label'] = 1
fake_df['label'] = 0

# Merge both datasets and then shuffle
(random_state to ensure reproducible
results).
df = pd.concat([true_df, fake_df],
axis=0).sample(frac=1,
random_state=42).reset_index(drop=True)

# Merge Title and Text columns together.
df['full_text'] = df['title'] + " " +
df['text']

# 2. NLP Preprocessing Pipeline
def clean_text(text):
    text = re.sub(r'https?://\S+|www\.\S+',
'', text)
    text = re.sub(r'<.*?>', '', text)
    text = re.sub(r'^\w\s', '', text)
    text = text.lower()
    words = text.split()
    words = [word for word in words if word
not in stop_words]
    return ' '.join(words)

# Apply the cleaning function
df['clean_text'] =
df['full_text'].apply(clean_text)

# 3. Train/Test Split (80/20)
X_train, X_test, y_train, y_test =
train_test_split(df['clean_text'],
df['label'], test_size=0.20,
random_state=42)

# 4. TF-IDF Vectorization
tfidf_vectorizer =
TfidfVectorizer(stop_words='english',
max_df=0.7, ngram_range=(1, 2))
tfidf_train =
tfidf_vectorizer.fit_transform(X_train)
```

```
tfidf_test =  
tfidf_vectorizer.transform(X_test)  
  
# 5. Initialize and Train the PAC Model  
pac =  
PassiveAggressiveClassifier(max_iter=50,  
random_state=42)  
pac.fit(tfidf_train, y_train)  
  
# 6. Predict and Evaluate  
y_pred = pac.predict(tfidf_test)  
  
# Calculate Metrics  
accuracy = accuracy_score(y_test, y_pred)  
conf_matrix = confusion_matrix(y_test,  
y_pred)  
  
print(f"Accuracy: {accuracy * 100:.2f}%")  
print("Confusion Matrix:\n", conf_matrix)  
print("Classification Report:\n",  
classification_report(y_test, y_pred))
```