



Automating fake news detection system using multi-level voting model

Sawinder Kaur¹ · Parteek Kumar² · Ponnurangam Kumaraguru³

Published online: 2 November 2019
© Springer-Verlag GmbH Germany, part of Springer Nature 2019

Abstract

The issues of online fake news have attained an increasing eminence in the diffusion of shaping news stories online. Misleading or unreliable information in the form of videos, posts, articles, URLs is extensively disseminated through popular social media platforms such as Facebook and Twitter. As a result, editors and journalists are in need of new tools that can help them to pace up the verification process for the content that has been originated from social media. Motivated by the need for automated detection of fake news, the goal is to find out which classification model identifies phony features accurately using three feature extraction techniques, Term Frequency–Inverse Document Frequency (TF–IDF), Count–Vectorizer (CV) and Hashing–Vectorizer (HV). Also, in this paper, a novel multi-level voting ensemble model is proposed. The proposed system has been tested on three datasets using twelve classifiers. These ML classifiers are combined based on their false prediction ratio. It has been observed that the Passive Aggressive, Logistic Regression and Linear Support Vector Classifier (LinearSVC) individually perform best using TF-IDF, CV and HV feature extraction approaches, respectively, based on their performance metrics, whereas the proposed model outperforms the Passive Aggressive model by 0.8%, Logistic Regression model by 1.3%, LinearSVC model by 0.4% using TF-IDF, CV and HV, respectively. The proposed system can also be used to predict the fake content (textual form) from online social media websites.

Keywords Fake news articles · Count–Vectorizer · TF–IDF · Hashing–Vectorizer · Classifiers · Textual content · Machine learning models

1 Introduction

A growing interest related to fake news detection has attracted many researchers as fake information is circulated through online social media platforms such as Facebook and Twitter. The fake content is spreading at a faster pace to

gain popularity over social media, to distract people from the current critical issues. Most of the people believe that the information they receive from various social media sites is reliable and true, i.e., people are inherently truth-biased. Also, people easily trust and want to believe in what they actually interpret in their minds, i.e., confirmation-biased. In general, it has been analyzed that humans are unable to recognize deception effectively. Due to this, a serious and negative impact of fake articles can be seen on society and individuals leading to an imbalance of the news ecosystem. It was observed that during US president election (Conroy et al. 2015), most of the widely spread articles on social media were fake. Recently, a fake video related to Kerala battling with floods was viral on social media platform (Facebook) claiming that the Chief Minister of the state is forcing the Indian Army to stop conducting the rescue operations in flooded regions of Kerala. Also, during India's national election (2019), various WhatsApp groups (> 900,000) were created to disseminate the fake information regarding India's ruling party (<http://bit.ly/2miuv9j>). Most of the fake articles

Communicated by V. Loia.

✉ Sawinder Kaur
sawinderkaurvohra@gmail.com

Parteek Kumar
parteek.bhatia@thapar.edu

Ponnurangam Kumaraguru
pk@iiitd.ac.in

¹ Doctoral Research Lab-II, Computer Science and Engineering Department, TIET, Patiala, India

² Computer Science and Engineering Department, TIET, Patiala, India

³ Computer Science and Engineering Department, IIIT, Delhi, India

are created to confuse people and trigger their distrust. Such problems led researchers to look at some automated ways to access the ground truth values of fake text on the basis of the textual content posted in articles on social platforms.

Social media enables to maintain and develop relations with others. It helps the users to present themselves by creating their profiles, sharing information through photographs, images and text to link with other members (Canini et al. 2011). Some of the most popular social media (Ahmed et al. 2017) sites are Facebook, Twitter (Gupta et al. 2013a; Wang 2010; Benevenuto et al. 2010), Instagram (Sen et al. 2018), WhatsApp (Garimella and Tyson 2018; Caetano et al. 2018), LinkedIn, WeChat, Snapchat and Foursquare (Pontes et al. 2012b). With the popularity of social media sites (Dewan et al. 2013), level of usage to share content on online social media has increased. There are several reasons for the change in behavior for such kind of consumptions. The content shared on social media platforms requires less time and cost than on newspapers or traditional news media. It is easier to share content in the form of video, blogs, posts with friends or users. This gives the growth ease to the authors and publishers to publish their contents as articles on collaborative environments. There is 13% of the global increase in social media usage since 2017 (Garimella and Tyson 2018). Distribution and creation of news content in the form of posts (Weimer et al. 2007), blogs, articles, images, videos, etc., have been spreading through social media websites (Dewan et al. 2013). This increase in social media also gives rise in the spread of fake articles (Wei and Wan 2017) over the Internet.

1.1 Types of fake news

According to the literature, there are five types of fakes news. The first type can be seen in the form of *deliberate misinformation*, which is misleading information that is spread in a calculated way to deceive targeted users. Other forms of fake news can be *clickbait* (Chen et al. 2015; Shu et al. 2017) which grab the reader's attention with a purpose to make them click on the fake news seen on the Internet. Users who set up fake sites generate huge revenues and clicking on such websites results in bombarded ads. Articles (Wei and Wan 2017) from satirical sources like 'The Onion' often repeat and share the news as if the printed stories were true. *Parody or Satirical* (Rubin et al. 2016) articles use obscenity, absurdity and exaggeration to comment on current events and to unease the readers. *False headlines* are intentionally exaggerated to draw the reader's attention. In such headlines, the title of the articles may not match the context of stories, the headline can be read as one way and state something different as a fact. This type of fake news is untrue at worst and misleading at best. *Hoaxes* is another type of misinformation which deceives the reader deliberately by causing harm and material losses to the users.

1.2 Contribution

Researchers have analyzed that an automated system sometimes identifies fake news articles better than human do. The automated systems can play an important tool for identifying fake stories, articles, blogs, clicks which manipulate public opinion on social networking sites (Dewan and Kumaraguru 2017; Jain and Kumaraguru 2016). Taking into need for the development of such fake news detection system, in this paper we have identified news articles as fake or real by using supervised machine learning classifiers such as Naïve Bayes (NB), Decision Tree (DT), Support Vector Machine (SVM), Linear models, Neural Networks (NN) and Ensemble models. To get effective results, three different corpora (News Trends, Kaggle and Reuters) have been collected with similar characteristics. The feature extraction techniques, Term Frequency–Inverse Document Frequency (TF–IDF), Count–Vectorizer (CV), Hashing–Vectorizer (HV), are used to extract feature vectors from the textual content of articles. The effectiveness of a set of classifiers is observed when they predict the labels of the testing samples after learning the training data using these features extraction techniques.

Various supervised ML models are extensively used for categorization of textual data as fake or real, but such models were not able to obtain the results of an ideal classifier. So, a multi-level voting model has been proposed in this paper to build an ideal classifier with significant improvement than the previously existing models. Our contribution to this paper is as follows:

- Statistical analysis of collected datasets (News Trends, Kaggle and Reuters) with negative and positive instances has been performed.
- Twelve ML models are evaluated using TF-IDF, CV, HV feature extraction techniques to retrieve the best model based on performance metrics.
- A novel method is proposed to merge the ML models based on false prediction rate.
- Proposed an ideal multi-level Voting Classifier (three-level) to verify the effectiveness of the ensemble model.
- A comparative study is performed to show the effectiveness of our proposed model.

The performance of the multi-level voting system is analyzed using parameters like precision, recall, specificity, ROC curve, *F1*-score.

The remaining paper is organized as follows: Sect. 2 gives a brief overview of the related work done in the field of fake news classification. The problem statement is discussed in Sect. 3. Section 4 covers the methodology of the proposed system. Section 5 introduces the classifiers used for detecting the fake and real articles. The evaluation phase along with analysis done by all classifiers is presented in Sect. 6. The

proposed model is discussed in Sect. 7, and its performance is evaluated in Sect. 8. The comparison of the proposed model with existing work is discussed in Sect. 9. Section 10 concludes the paper along with its future works.

2 Related work

Many techniques have been developed recently for fake news detection. In this section, the closely related research work that has been done on detecting fake news on online social media is discussed.

Most of the social sites require energy and time to manually remove or filter spam. Markines et al. (2009) proposed six highlights (TagSpam, TagBlur, DomFp, NumAds, Plagiarism, ValidLinks) of tagging systems catching diverse properties of social spam (Markines et al. 2009). Utilizing the six proposed highlights, creators assessed different administered machine learning techniques to identify spam with precision over 98% with a false positive rate of 2%. The Weka tool used for the experiment gives the best accuracy when evaluated on the basis of AdaBoost classifier. To address the issue of detecting video promoters and spammers, Benevenuto et al. (2009) manually assembled test gathering of genuine YouTube clients and classified them as legitimates, spammers and promoters. Authors have investigated the feasibility for detecting spammers and promoters by applying a supervised classification algorithm (Benevenuto et al. 2009). The classifier used in this paper correctly identifies the majority of promoters correctly.

Qazvinian et al. (2011) explored three features such as content-based, network-based and microblog-specific memes to correctly identify the rumors (Qazvinian et al. 2011). Such features are also used to identify disinformers or users who endorse a rumor and further tries to spread it. For the experiment, authors collected 10,000 manually annotated tweets from twitter and achieved 0.95 in mean average precision (MAP). Rubin et al. (2016) proposed a satire detection model with Support Vector Machine (SVM)-based algorithm across 4 domains, such as science, business, soft news and civics (Rubin et al. 2016). To verify the sources of news articles, authors have discussed various legitimate and satirical news websites. In this paper, five features together are chosen to predict the best predicting feature combination with 90% precision and 84% recall to identify satirical news which can help to minimize deception impact of satire.

Analysis for the real event such as Boston Marathon Blasts is done by Gupta et al. (2013a). During the event, it was observed that a lot of fake and malicious profiles originated on Twitter. Results showed that 29% of the content originated during the Boston Blasts (Gupta et al. 2013a) was viral, whereas 51% was generic opinions and comments. Authors

identified six thousand profiles which were immediately created after the blasts occurred and were suspended by Twitter. Tabloids are often used for sensationalization, exaggeration, producing misleading and low-quality news content. A new form of tabloidization has emerged known as clickbaiting. There exists both non-textual and textual clickbaiting, which is surveyed by Chen et al. (2015), who proposed a hybrid approach (Chen et al. 2015) for automatic detection of clickbait.

Rubin et al. (2015) utilized vector space modeling (VSM) and rhetorical structure theory (RST) to analyze misleading and truthful news. RST catches the coherence of story in terms of functional relations among the useful text units and also describe the hierarchical structure for each article/news. VSM is used for identifying the relations among rhetorical structure (Rubin et al. 2015), i.e., each article content can be depicted as vectors in high-dimensional space.

Researchers have used different techniques to identify and review the fake content. One of the best and common feature extraction method is Bag of Words. This comprises a group of words retrieved from the textual content, from where n-gram (Ahmed et al. 2017) features can be extracted. The second most important feature which is similar to the Bag of Words approach is Term Frequency (TF) which is related to frequency of the words. Conroy et al. (2015) proposed a hybrid approach which combines both machine learning and linguistic cues with network-based behavioral data (Conroy et al. 2015). The hybrid approach follows both n-gram and Bag of Words techniques to represent data. Ahmed et al. (2017) proposed a fake news detection system which uses n-gram (Magdy and Wanas 2010) analysis and Term Frequency–Inverse Document Frequency (TF–IDF) as a feature extraction technique (Ahmed et al. 2017). In this paper, six classifiers of machine learning are used and two different feature extraction techniques are used for comparison and investigation. Volkova et al. (2017) built a predictive model to manage 130K news posts as verified or malicious. Authors have classified four subtypes of suspicious news such as propaganda, clickbait, hoaxes and satire (Volkova et al. 2017).

Chhabra et al. (2011) had put forward a URL static feature-based detection method for detecting malicious websites with accurate results. The author has focussed on external features such as IP addresses. Further, a vector construction VSM (Chhabra et al. 2011) is chosen as the URL vector model. The dataset taken in this paper consists of malicious URLs which was downloaded from the phishing platform named as ‘Phishtank’ (Aggarwal et al. 2012).

In our digital world, fake news is disseminating and impacting millions of users on social media platforms every day. It has really become difficult to separate truth from fiction. With the help of machine learning models, it is possible to detect spam emails at an early stage with the help of spam filters. ML classifiers help to solve the real-world problems.

Table 1 Comparative analysis of research studies

Authors	Proposed approach	Model	Dataset	Features
Markines et al. (2009)	Analyzed distinct six features for detecting social spammers using machine learning	SVM, AdaBoost	Spam posts, tags	TagSpam, TagBlur, DomFp, NumAds, Plagiarism, ValidLinks
Benevenuto et al. (2009)	A video response crawler is proposed to identify spammers in online video social network	SVM	Real YouTube user information	Video attributes, individual characteristics of user behavior, social relation between users via video response interactions
Qazvinian et al. (2011)	Identified tweets in which rumor is endorsed	Naïve Bayes	Tweets	Content-based, network-based, Twitter specific memes
Chhabra et al. (2011)	Using URLs static features, a method is developed to detect malicious websites	Naïve Bayes, Logistic Regression, DT, SVM-RBF, SVM-Linear, SVM-Sigmoid	Malicious URL dataset from 'Phishtank'	Grammar, Lexical, Vectors and Static
Gupta et al. (2013a)	Analysis of Twitter content during Boston Marathon	Logistic Regression	Tweets and corresponding user information	Topic engagement, Global engagement, Social reputation, Likability, Credibility
Chen et al. (2015)	Analyzed coherence relations between deceptive and truthful news	VSM	News samples from NPR's 'Bluff the Listener'	Discourse
Rubin et al. (2015)	A hybrid approach is proposed combining linguistic and network-based behavior data	Linguistic, Network models	Simple text sentences	Bag of Words, n-gram
Conroy et al. (2015)	A satire detection model is developed	SVM	US and Canadian national newspapers	Absurdity, Humor, Grammar, Negative affect, Punctuation
Ahmed et al. (2017)	Developed n-gram-based classifier to differentiate between fake ad real articles	LinearSVM	News articles	TF-IDF
Caetano et al. (2018)	A predictive model was built to predict 4 subtypes of suspicious news: satire, hoaxes, clickbait and propaganda	Linguistic models	News posts	TF-IDF, Doc2Vec
Proposed system	Using textual data of articles, an efficient multi-level voting model is developed to detect fake articles	SGD, PA, MultinomialNB, Gradient Boosting, DT, AdaBoost	News articles	TF-IDF, Count-Vectorizer, Hashing-Vectorizer

Also, ML has made easier for the users of e-commerce business as it helps to identify the hidden pattern, groups the similar products into a cluster and displays the result to end user which enables a product-based recommendation system. It also helps to solve the problem of unfair recommendations (D'Angelo et al. 2019).

Comparative analysis of related work done in the field of fake news detection and the proposed system presented in this paper is shown in Table 1.

It has been analyzed that the research work done in the field of fake news detection is mainly restricted to SVM and Naïve Bayes classifiers using only n-gram (Magdy and Wanas 2010) and TF-IDF features extraction approaches. No work has been done on Multi-Layer Perceptron (MLP), Long Short-term Memory (LSTM) (LeCun et al. 2015) models

and hashing-based extraction approach which also accounts for better efficiency. Also, the existing fake news detection (Ruchansky et al. 2017) models are built using supervised machine learning algorithms, whereas manual hand-crafted feature extraction is more time-consuming and inefficient method to achieve the best accuracy. The existing techniques studied so far provide a direction to be followed further for quantitative and qualitative research.

3 Problem statement

To address the issue of fake news generation and dissemination through various online social platforms, an appropriate feature extraction technique is chosen to improve the effi-

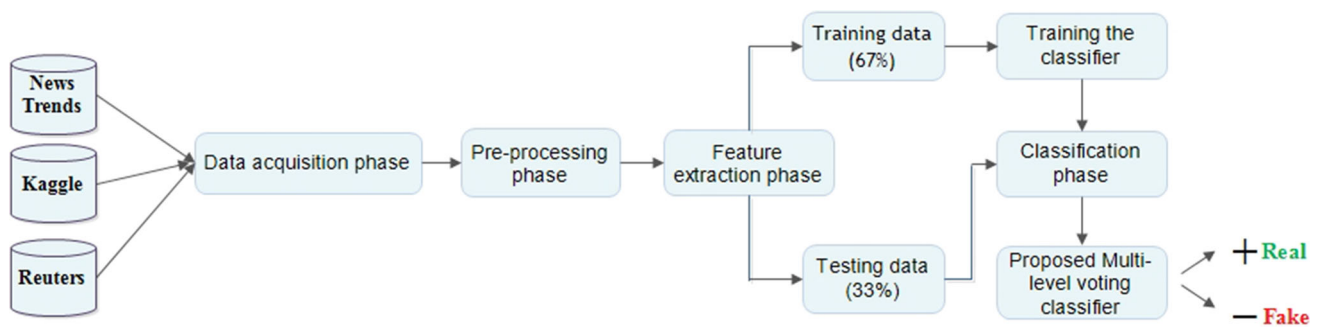


Fig. 1 Architecture of the proposed automatic fake news detection system

ciency of the existing ML classifiers. A novel multi-level voting ensemble model will be proposed to develop an efficient fake news detection system. Mathematically, the problem statement can be represented as follows: to identify $S = \{\text{fake}, \text{real}\}$ for a document D where $D = \{t_1, t_2, \dots, t_n\}$ and t_i represents the text in a news article a_i chosen from corpus with series of engagements that is composed of title, body and label of the article as $e_{ijk} = (t_i, b_j, l_k)$. The task is to evaluate and analyze the best feature extraction method f_m where $m = \{\text{TF-IDF}, \text{CV}, \text{HV}\}$ using machine learning classifier to compute high efficiency in our proposed system. The approach followed in this paper will be discussed in the next section.

4 Methodology

The architecture of the proposed fake news article detection system (Ahmed and Abulaish 2012) is shown in Fig. 1. To train the system, three corpora have been collected from three different sources by downloading the datasets from News Trends, Kaggle and Reuters websites. In the pre-processing phase, stop words and duplicate text from news articles are removed. The missing values, i.e., not available (NA) values, are collected and cleaned in the next step. The data retrieved are then split into two parts, training (0.67) and testing (0.33) sets. The feature extraction phase is then carried out to retrieve meaningful features from the textual data. In this phase, the features are extracted from the articles. Three feature extraction techniques such as Term Frequency-Inverse Document Frequency (assigns weights according to the importance of the terms in the document), Count-Vectorizer (counts the frequency of the terms in a document) and Hashing-Vectorizer (follows the hashing trick) have been applied. The features retrieved are then fed to the classification algorithm chosen in next phase. The various ML models such as MultinomialNB, Passive Aggressive, Stochastic Gradient Descent, Logistic Regression, Support Vector Classifier, Nu-Support Vector Classifier, Multi-Layer Perceptron, Linear Support Vector Classifier, AdaBoost, Gra-

dient Boosting, Decision Tree, Voting Classifiers (Sirajudeen et al. 2017) are chosen to learn and identify the patterns and outcomes from them. The models are then evaluated based on performance metrics to achieve an efficient classifier. Based on the analysis done, the models are integrated to propose a multi-level voting model to achieve high efficiency and then is compared with the existing work (Gao et al. 2010) as discussed in Sect. 9. The detailed working of each phase that has been implemented in python framework is discussed below.

4.1 Data collection

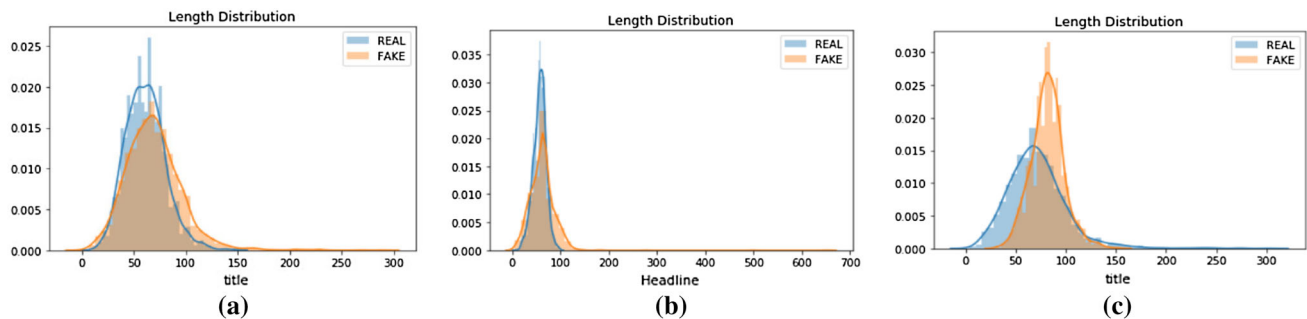
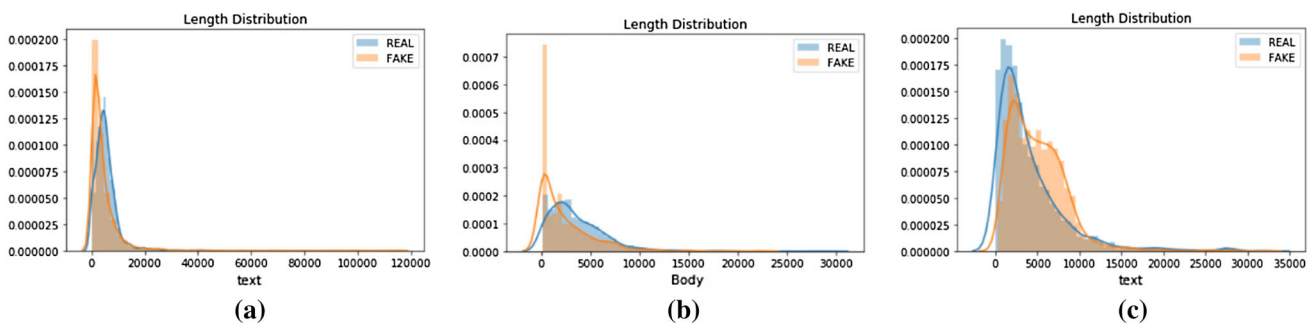
There are many sources of fake article generation such as Facebook (Dewan and Kumaraguru 2015) and Twitter (Aggarwal et al. 2018; Gupta and Kumaraguru 2012b, a; Gupta et al. 2013b), which are used as a trading platform to disseminate fake news. We have used News Trends (<https://bit.ly/2zVRLxK>), Kaggle (<https://bit.ly/2Ex5VsX>) and Reuters (<https://bit.ly/2BmqBQE>) dataset with similar attributes such as headlines, body, publisher name of the article, published date, categorical and missing values. The News Trends, Kaggle and Reuters corpus consist of 7,795; 4,048 and 21578 news articles labeled as fake and real news, respectively. Overall statistics of our three collected datasets are discussed in Table 2. To analyze the length distribution of the titles for fake and real articles, Fig. 2 is visualized, where the X-axis labeled as 'title' represents the number of terms used in news article titles or headlines, whereas Y-axis represents the corresponding number of articles having the same length distribution. A conclusion can be drawn from the mean distribution by visualizing Fig. 2 that the length of the titles or headlines of fake news articles is often longer than real news articles.

To further analyze the length distribution of the body content in articles, Fig. 3 is observed, where the X-axis labeled as 'text' represents the number of terms used in news article texts, whereas Y-axis represents the corresponding number of articles having the same length distribution.

It has been analyzed that the length of body/text of real news articles are often higher than that of fake news articles

Table 2 Statistics of collected corpora

Corpus	Total article	Cleaned articles	Real articles	Fake articles	Published year
News Trends	7795	6335	3171	3164	2017
Kaggle	4048	3983	1865	2118	2017
Reuters	21,578	19,969	9622	10,347	2004

**Fig. 2** Length distribution of headline for fake and real articles on **a** News Trends, **b** Kaggle and **c** Reuters corpora**Fig. 3** Length distribution of text for fake and real articles on **a** News Trends, **b** Kaggle and **c** Reuters corpora

as shown in Fig. 3a, b, but for Reuters corpus the length of real news articles is more than fake news articles as seen in Fig. 3c.

The statistics of the mean distribution for Figs. 2 and 3 are compared in Table 3. In general, a conclusion can be drawn after analyzing different datasets that the headlines of fake articles are longer in length and have shorter body content than real articles published on social networking (Pontes et al. 2012a) sites.

4.2 Pre-processing

In the pre-processing phase, non-semantic words such as prepositions, conjunctions, and pronouns also known as stop words are removed from the textual document as they provide very little or no information about fake content of an article. The redundant data in the form of textual strings are removed from the document using a regular expression (*regex*) in the next step as shown in Fig. 4. The *regex* and *pandas* library has been used to perform the pre-processing task (Batchelor 2017). *Regex* library has been used in python to define

a search pattern using a sequence of characters, whereas *dropna* method from *pandas* is used for cleaning the missing values in python DataFrame. We have used *random_state* function to select the entries from the dataset which is further used to split training and testing data points as it is used to split the data randomly.

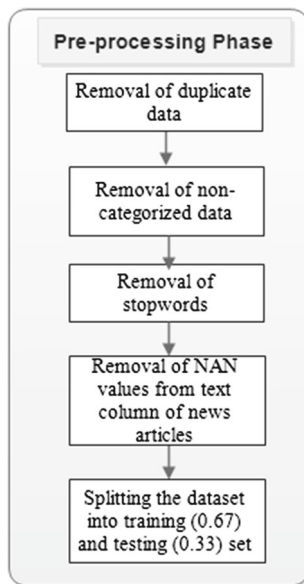
To avoid overfitting, three standard splits (70:30, 67:33 and 60:40) were used to perform the experiment. When the first standard split (70:30) was performed, it was observed that the data point dealt with an issue of underfitting. During the second split (60:40), overfitting of data was analyzed, whereas the third split (67:33) gave the best predicted line covering the majority of data points in the graph, so a standard split of 67:33 was chosen. The training data are then fed to feature generation phase as discussed in the next section.

4.3 Feature generation

In order to extract numerical features from a textual document, tokenization, counting and normalization are done. During tokenization, each word is given a unique integer

Table 3 Mean distribution of labeled articles for News Trends, Kaggle and Reuters corpora

Corpus	Mean distribution of title labeled as fake	Mean distribution of title labeled as real	Mean distribution of text labeled as fake	Mean distribution of text labeled as real
News Trends	69.18	61.38	4121.04	5292.16
Kaggle	62.47	57.32	2380.82	3544.84
Reuters	81.79	71.29	5156.08	4540.26

**Fig. 4** Steps performed during the pre-processing phase

ID, following which occurrence of tokens is counted and then normalization of such tokens takes place. The whole process of converting the textual document into numerical feature vector is called as vectorization. Together, this strategy (tokenization, counting, normalization) is called as ‘Bag of n -grams’ or ‘Bag of Words’ where n represents the continuous sequence of terms (Alahmadi et al. 2013). The three feature extraction techniques as shown in Fig. 5 for retrieving features from the textual content of an article are Count-Vectorizer, TF-IDF and Hashing-Vectorizer which use *CountVectorizer*, *TfidfVectorizer* and *HashingVectorizer* classes from *feature_extraction* library of python, respectively.

4.3.1 Count-Vectorizer (CV)

Count-Vectorizer is absolutely based on count of the word occurrences in the document. In Count-Vectorizer technique, both counting the occurrences of tokens and tokenization process are performed. Count-Vectorizer has many parameters for refining the type of features. One can build features using any of the three parameters, unigram ($min_df = 1$), bigram ($min_df = 2$) and trigrams ($min_df = 3$). We have used

min_df as 1 for our experiment. Here, each vector (term) in a document represents the individual feature name and its occurrence is depicted through a matrix to make it easier to understand as shown in Table 4. It has been observed from the above table, after pre-processing phase, the terms retrieved from the documents are represented as vectors on top of the sparse matrix and the frequency of terms in particular document is represented through count occurrences. Tag clouds of top 30 features retrieved after executing CV method are shown in Fig. 6. It was observed that words like *corruption*, *attacking*, *islamic*, *obama*, *losing*, *com* are seen under fake tag clouds.

CV also counts the number of words occurring more frequently in the document, which may overshadow the words which occur less frequently but may have more importance to the document feature. This limitation of CV can be handled by using TF-IDF feature extraction technique as explained below.

4.3.2 Term Frequency–Inverse Document Frequency (TF-IDF)

TF-IDF is a weighing matrix which is used to measure the importance of a term (count + weight) to a document in a dataset. The tokens retrieved from the textual data using both TF-IDF and CV techniques are same, but weights assigned to the tokens of both techniques are different. TF-IDF is composed of two metrics, named as term frequency (tf) and inverse document frequency (idf). TF-IDF is represented by Eq. (1).

$$tfidf = tf(t, d) \times idf(t, d) \quad (1)$$

Here, Term Frequency is denoted as tf and is calculated from the count (c), term (t) in document (d) and represented as $tf(t, d) = c_{td}$. The frequency of occurrence of words to a binary feature is converted by using 1 (present in document) and 0 (not present in document). The frequencies can be normalized using average and logarithms. The inverse document frequency (idf) for a word w in document text (t) as computed by Eq. (2).

$$idf(t, d) = 1 + \log \frac{T}{(1 + df(t))} \quad (2)$$

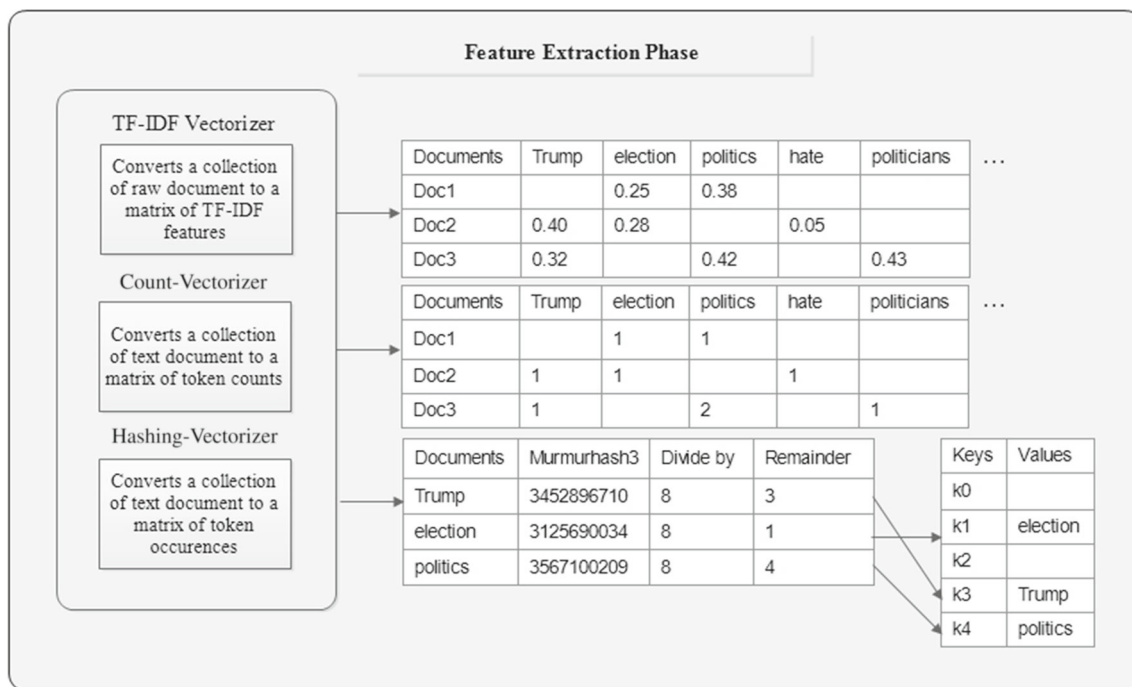


Fig. 5 Conversion of textual content into numerical vectors through TF-IDF, Count-Vectorizer and Hashing-Vectorizer feature extraction techniques

Table 4 Sparse matrix representation using a Count-Vectorizer feature extraction technique

Document	Narendra	Elections	Vote	Politics	Punjab	BJP	Candidate
Doc1	0	1	2	0	0	2	1
Doc2	4	0	1	0	0	1	0
Doc3	1	3	2	0	1	0	1

Here, T represents the total document count in our corpus and $df(t)$ represents the count of the number of documents where the term t is present. The product of two measures will help to compute $tfidf$. Euclidean's normalized form is used to calculate the final TF-IDF metric as given by Eq. (3).

$$tfidf = \frac{tfidf}{||tfidf||} \quad (3)$$

Here, $||tfidf||$ is the Euclidean norm. Tag clouds of top 30 features for fake and real articles using TF-IDF feature extracting technique are shown in Fig. 7. It was observed that words like *www*, *danger*, *sorry*, *wars*, *islamic* are most common under fake news articles. The difference between Count-Vectorizer and TF-IDF approach is that the tokens retrieved from textual data are same, but both have different weights assigned to each token being extracted.

4.3.3 Hashing-Vectorizer (HV)

It is a memory-efficient technique. Unlike previous two techniques, tokens are stored as a string and here hashing trick

is applied to encode the features as numerical indexes. Let us discuss the concept of HV using an example shown in Fig. 8. When data are entered, the hashed attributes of data are retrieved. The hashed terms like Trump, election and politics are extracted from the document. In the next step, the hashing trick is applied on the hashed attributes, where a *Murmurhash3* function is applied on the hashed terms to generate a random number. Further, the assigned random numbers are divided by 8 and are stored in different keys such as k_2, k_3, k_4 based on the remainders retrieved after applying the *murmurhash3* function which is used for hash-based lookup. There is a possibility of collision when data have equal hashed attributes.

Let us suppose in our example document, we have Trump and politics words as an important keys which are seen more than once, thus causing collisions at k_3, k_4 positions. The collided values are then occupied by other vacant positions in a set of documents. Such collision processing is dealt with parallel processing. The process is conceptually explained in Fig. 9. These six terms are assigned six keys as shown in Fig. 9 and are entered in the hash table. The hash values of keys k_1, k_3, k_6 are same, i.e., *Trump*; k_4 and k_5 are same,

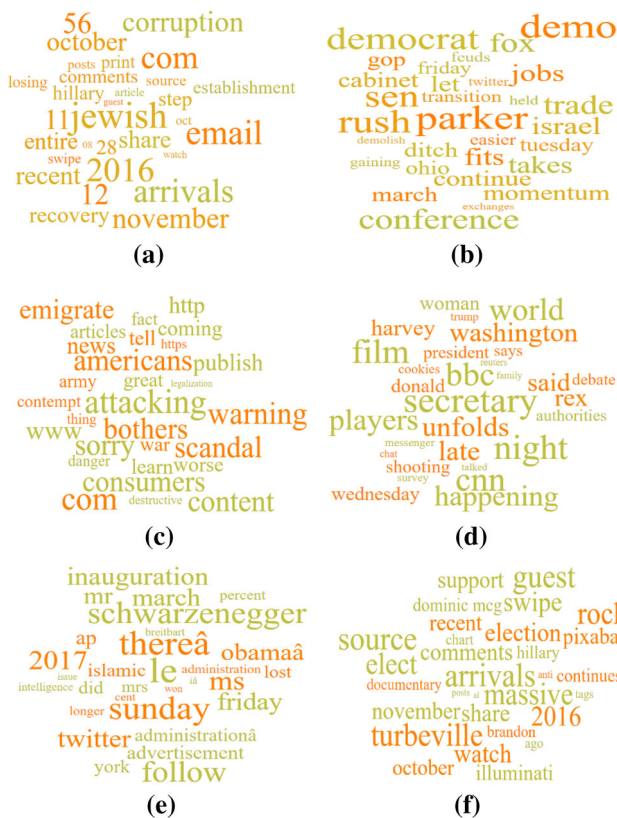


Fig. 6 Top 30 **a** News Trends fake, **b** News Trends real, **c** Kaggle fake, **d** Kaggle real, **e** Reuters fake and **f** Reuters real-word clouds using Count-Vectorizer feature extraction technique

i.e., politics, but rest have different values. Due to collision, k_1 , k_3 and k_6 cannot be placed in same set (S_1). To enable parallel processing, k_1 , k_3 and k_6 are placed in different sets. No two same hash values can be placed in a single set. Different keys like k_1 , k_2 and k_4 can be placed in the same set (S_1) as these keys have different values. Values at keys, k_3 and k_6 , are same, so cannot be processed parallelly; therefore, are processed in different sets. The values in S_1 , S_2 , S_3 are organized into vectors (numerical features) and can be processed using vector operations.

The drawback of HV is that there is no way to get the feature names from feature indices, i.e., the inverse transform cannot be used to compute the most important feature names through Hashing-Vectorizer, unlike rest two methods.

5 Classification algorithms

The processed dataset retrieved after pre-processing and feature extraction phase is then fed to the classification phase for the identification of fake news article. In this paper, six machine learning techniques, i.e., Naïve Bayes (MultinomialNB), Support Vector Machine [Support Vector Classifier



Fig. 7 Top 30 **a** News Trends fake, **b** News Trends real, **c** Kaggle fake, **d** Kaggle real, **e** Reuters fake and **f** Reuters real-word clouds using TF-IDF feature extraction technique

Documents	Murmurhash3	Divide by	Remainder	Keys	Values
Triump	3452896710	8	3	k0	-
election	3125690034	8	2	k1	-
politics	3567100209	8	4	k2	election
				k3	Triump
				k4	politics

Fig. 8 Hashing trick on textual data using Murmurhash3 function to get values in a specific range

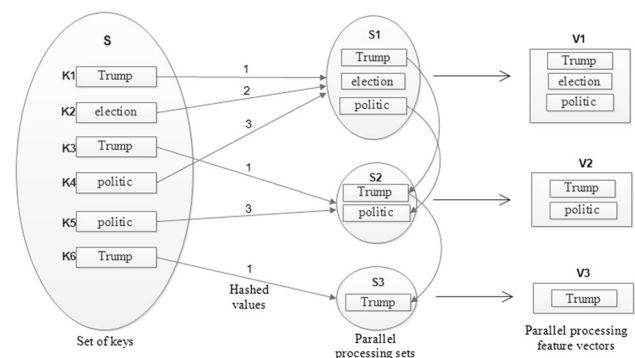


Fig. 9 Entry of redundant data into a hash table during parallel processing

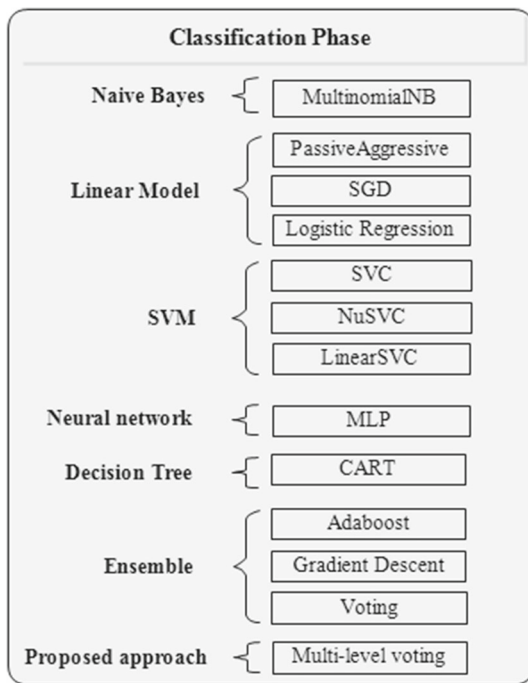


Fig. 10 Classification phase

(SVC), NuSVC, LinearSVC], Decision Tree (CART), Linear [Passive Aggressive (PA), Stochastic Gradient Descent (SGD), Logistic Regression (LR)], Neural Network [Multi-Layer Perceptron (MLP)] and Ensemble models (AdaBoost, Gradient Boosting, Voting) have been applied as shown in Fig. 10.

Classifier functions help to map the input feature vectors $f \in F$ to output labels $l \in \{1, 2, 3, \dots, n\}$, where F is the feature space. The feature space is represented as $F = \{\text{Fake}, \text{Real}\}^R$, where R is the real number. Our aim is to learn the classifier function from labeled training data.

5.1 Naïve Bayes (NB)

It is a type of probability classifier. It works on Bayes theorem and handles both categorical and continuous variables. NB assumes that every pair of features with labeled value is independent of each other. Given a collection of D documents from news articles, $D_i = \{d_1, d_2, \dots, d_n\}$, where each document consists of T terms such as $D_i = \{t_1, t_2, \dots, t_m\}$. Then, the probability of D_i occurrence in class label C_l is given by Eq. (4).

$$P(C_l|D_i) = P(C_l) \prod_{n=1}^m P(d_n|C_l) \quad (4)$$

Here, the conditional probability of term t_m present in a document of class label C_l and the prior probability of document occurring in class label C_l is denoted by $P(C_l)$.

Multinomial Naïve Bayes (MultinomialNB) is a type of Naïve Bayes algorithm used for text classification. Data used in text classification to apply MultinomialNB can be represented as *TF-IDF* vectors, hashing vectors and count vectors. The feature vectors $V_f = (V_{f1}, V_{f2}, \dots, V_{fn})$ are parameterized for each class C_n in the distribution, where n represents the feature numbers. The likelihood of observing V_f is given by Eq. (5).

$$P(C_n|V_f) = \frac{X_i!}{\prod_{n=1}^p f_n!} \prod_{n=1}^p V_n f_n \quad (5)$$

Here, f_n is the number of times the n th feature has occurred in the document, X_i is the number of draws taken from the bag of features. $V_n^{f_n}$ and $f_n!$ are computed from training data.

5.2 Linear model

The linear model helps to classify the group by making linear combinations of feature vectors. Linear classifiers work well with many features, but it works best for document (features are extracted from the text) classification. If \mathbf{v} is the input feature vector to the classifier, then the resultant score is given by Eq. (6).

$$s = f(\mathbf{v}\mathbf{w}) = f\left(\sum_i w_i v_i\right) \quad (6)$$

here \mathbf{w} is the weight of a feature vector and the function f gives the desired output of two vectors. The three linear models used in this paper are Passive Aggressive (PA), Stochastic Gradient Descent (SGD) and Logistic Regression (LR) classifiers. The PA algorithm has similar behavior with perceptron classifier in terms of learning rate but has dissimilar behavior in terms of regularization parameter. The PA classifier is equivalent to PA-I (Dewan and Kumaraguru 2015) when the loss parameter is *hinge* and PA-II (Dewan and Kumaraguru 2015) when the loss parameter is *squared_hinge*. Second linear model used in this paper is SGD. The SGD model is updated with the decreasing learning rate after each sample interval and by default the loss parameter used in this paper is *hinge*. The classifier also allows minibatch learning. The LR model can be used for both multi- and binary problem classification. Any other input format apart from *float64* gets converted in this classifier. All three linear models discussed above can take both sparse and dense matrix as their input.

5.3 Support vector machine (SVM)

SVM works on structural risk minimization principle. It is defined by a 'best' separating hyperplane and is also called as

a discriminative classifier. Through the SVM model, feature vectors retrieved from text document of news articles are represented as points in feature space. Then, the feature vectors are mapped in such a way that a wide gap is visible to perform linear classification. In our dataset, the feature vectors are marked by making two categories, $C = \{\text{Fake}, \text{Real}\}$, and then the training classifier builds a model which assigns new feature vectors to both defined categories.

SVM classes such as Linear Support Vector Classifier (LinearSVC), Nu-Support Vector Classifier (NuSVC), Support Vector Classifier (SVC) are used for performing classification on the dataset. NuSVC and SVC are almost similar but use slightly different sets of parameters and their mathematical formulations also vary, whereas LinearSVC is another type of Support Vector Classification (SVC) and uses the case of linear kernel. All three classes take input in the form of two arrays with array X having two-dimensional size [sample_number, feature_vectors] to handle training data and array Y having two-dimensional size [category_label, sample_number]. The decision function that is same for both SVC and NuSVC is given by Eq. (7).

$$\text{sgn} \left(\sum_{f=1}^n y_f \alpha_f K(V_f, V) + \mu \right) \quad (7)$$

where V_f are the training feature vectors, $f = 1, 2, \dots, n$ in two categories. $K(V_f, V)$ is the kernel and $y_f \alpha_f$ is the dual coefficient parameter which holds support vectors and an independent intercept term μ . The only difference between SVC ($C = [0, \infty]$) and NuSVC ($C = [0, 1]$) is seen from parameter C which is the penalty parameter of the error term. The class LinearSVC supports both sparse and dense input and is implemented in terms of liblinear, so is more flexible in terms of the loss function and penalties, and scales better to large testing samples.

5.4 Neural network (NN)

Neural networks are composed of highly interconnected neurons to solve specific problem parallelly. In this paper, Multi-Layer Perceptron (MLP) is implemented on our collected dataset. The classifier can be trained on either regression or classification dataset. The feature vectors $V_f = \{v_1, v_2, \dots, v_n\}$ are retrieved after feature extraction phase, and through training dataset the classifier learns a given function in Eq. (8).

$$f(n) : Ri \rightarrow Ro \quad (8)$$

where i is the dimensions for input and o is the dimensions for an output. In MLP, there can be one or more nonlinear layer (hidden layer) between the input and output layer. The input

layer is made up of neurons, where each neuron represents the input feature which is fed into the hidden layer. Then, the hidden layer computes the weighted summation $w_1 v_1 + w_2 v_2 + \dots + w_i v_i$, followed by function $f(n)$. The output value is given by the last hidden layer and is received by the output layer.

5.5 Decision tree (DT)

The DT classifiers can be used for both regression and classification. The classifier predicts the target variable by learning the feature data and dividing the area into sub-regions. Based on two criteria, multiple features are divided: one is a measure of impurity, and other is information gain. In our dataset, 'gini' is the chosen impurity measure to calculate the highest information gain at each node for dividing the DT. In case of document data, the conditions depend upon the particular term in a text document of the news article. The data are divided repetitively until the leaf node cannot be further divided acquiring least information on them. The majority count of labels in the leaf node are used for classifying the textual data.

5.6 Ensemble methods

Such methods help to build a learning algorithm by combining the estimators to get robust classifier over single classifier. The boosting methods implemented on our dataset are Gradient Descent Boosting (GDB) and AdaBoost classifier for binary classification. The AdaBoost classifier assigns more weight to feature vectors which is difficult to handle and less weight to the features which can be easily handled. This process is repeated until the classifier correctly classifies the training data. The GDB model works on three elements such as weak learner, loss function and additive model as discussed in Shu et al. (2017).

The other type of classifier that can be useful for balancing individual weakness of ML models is Voting Classifier. In this paper, the Voting Classifier has predicted the categories based on 'hard' voting which classifies the sample based on majority class label. To evaluate the classifiers discussed above, the performance metrics are defined and corresponding experimental results are discussed in the next section.

6 Evaluation phase

The performance measures for binary classifiers applied in this paper have been evaluated with the help of a confusion matrix defined by four cells as shown in Table 5, where

Table 5 A confusion matrix representation

Actual↓ Predicted→	Fake	Real
Fake	TP (a)	FP (b)
Real	FN (c)	TN (d)

- cell ‘a’ counts the predicted document as ‘Fake’ when actually it is ‘Fake,’ known as true positive (TP) rate.
- cell ‘b’ counts the predicted document as ‘Real’ when actually it is ‘Fake,’ known as false positive (FP) rate.
- cell ‘c’ counts the predicted document as ‘Fake’ when actually it is ‘Real,’ known as false negative (FN) rate.
- cell ‘d’ counts the predicted document as ‘Real’ when actually it is ‘Real,’ known as true negative (TN) rate.

The conventional performance measure has been evaluated from the above confusion matrix cells. The measures computed from the matrix are precision as represented by Eq. (9), recall by Eq. (10), specificity by Eq. (11), accuracy by Eq. (12), error by Eq. (13), $F1$ -score by Eq. (14) as shown below.

$$\text{Precision (Pr)} = \frac{a}{a + b} \quad (9)$$

$$\text{Recall (Re)} = \frac{a}{a + c} \quad (10)$$

$$\text{Specificity (Sp)} = \frac{d}{d + b} \quad (11)$$

$$\text{Accuracy (Acc)} = \frac{a + d}{n}, \text{ where } n = a + b + c + d > 0 \quad (12)$$

$$\text{Error (Err)} = \frac{b + c}{n}, \text{ where } n = a + b + c + d > 0 \quad (13)$$

$$F1\text{-score (F1)} = \frac{2 \times \text{Pr} \times \text{Re}}{\text{Re} + \text{Pr}} \quad (14)$$

The predictions made by the classification models are evaluated in this phase based on their performance metrics. The most intuitive performance measure is accuracy, which helps

to predict the best model. Various machine learning models used in experiment are MultinomialNB (C1), Passive Aggressive (C2), Stochastic Gradient Descent (C3), Logistic Regression (C4), SVC (C5), NuSVC (C6), LinearSVC (C7), Multi-Layer Perceptron (C8), Decision Tree (C9), AdaBoost (C10), Gradient Descent (C11), Voting (C12) and multi-level voting (C13) classifiers. To evaluate these models, a comparative analysis is shown in Fig. 11. The experiment is performed on three (News Trends, Kaggle and Reuters) different corpus. In this paper, TF-IDF, CV and HV feature extraction techniques are used to extract the feature vectors from the documents of the chosen corpus. In Table 6, the accuracy measure of various ML classifiers is compared.

The best accuracy retrieved by top 3 models in all three corpora is Linear Support Vector Classifier (LSVC), Passive Aggressive (PA) and Logistic Regression (LR) classifiers.

Other parameters used to evaluate the performance measure of classifiers used in this paper are precision, recall and $F1$ -score. The precision metric helps to calculate the proportion of news article that are predicted fake and actually also belongs to the fake news article category. The comparative analysis of precision metric is shown in Table 7.

The recall metric helps to calculate the proportion of news article that are predicted fake but actually belongs to both fake and real articles. The comparative analysis of recall metric is shown in Table 8. The specificity metric helps to calculate the proportion of news article that are correctly predicted as a real news article known not to be fake. Specificity is measured as inverse of recall metric. The comparative analysis of specificity metric is shown in Table 9. Accuracy is only measured as a strong metric when both false negative and false positive values are closer to each other else the metric is not considered as a good performance measure. To convey a balance between recall and precision, $F1$ -score performance metric is selected to retrieve the best model. $F1$ -score metric helps to take both false positives and false negatives into account. The comparative analysis of $F1$ -score metric is shown in Table 10.

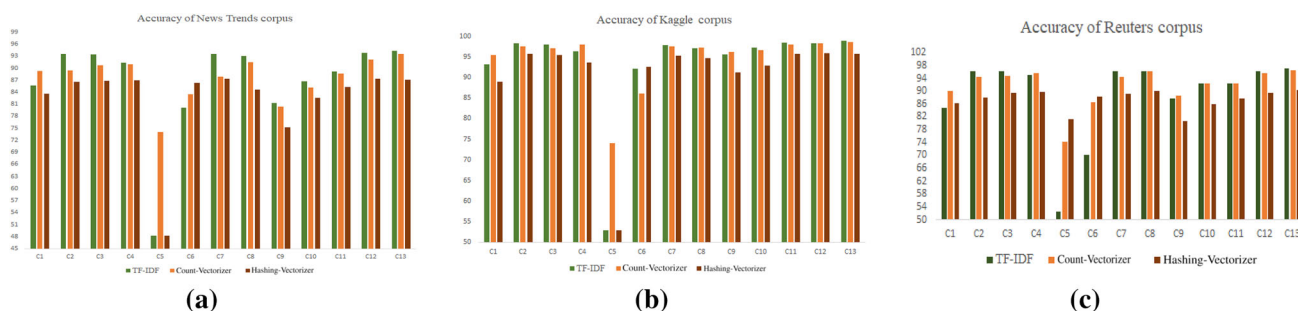


Fig. 11 Performance analysis using accuracy metric for **a** News Trends, **b** Kaggle and **c** Reuters dataset on the basis of TF-IDF, Hashing-Vectorizer and Count-Vectorizer feature extraction techniques

Table 6 Comparative analysis of accuracy measure using machine learning classifiers

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Multinomial Naïve Bayes	85.7	89.3	83.6	93.2	95.4	89	84.8	89.9	86.2
Passive Aggressive	93.5	89.4	86.6	98.3	97.6	95.7	96.2	94.3	88
Stochastic Gradient Descent	93.4	90.7	86.8	98	97.1	95.5	96.2	94.8	89.4
Logistic Regression	91.4	91.0	87.0	96.4	98	93.6	94.9	95.7	89.6
Support Vector Classifier	48.2	74.1	48.2	52.9	74	52.9	52.5	74.3	81.3
NuSVC	80.1	83.5	86.3	92.2	86.1	92.6	70	86.6	88.3
LinearSVC	93.6	87.9	87.3	97.9	97.6	95.3	96.3	94.4	89.3
Multi-Layer Perceptron	93	91.5	84.7	97.1	97.3	94.7	96.2	96.2	90.1
Decision Tree	81.3	80.4	75.2	95.6	96.2	91.3	87.7	88.5	80.6
AdaBoost	86.7	85.1	82.6	97.3	96.6	92.9	92.5	92.4	85.9
Gradient Boosting	89.2	88.6	85.3	98.5	98.0	95.7	92.3	92.5	87.6
Voting Classifier	93.8	92.1	87.3	98.3	98.3	95.9	96.1	96.4	90.3

Bold indicates the best performance metric values among the other models

Table 7 Comparative analysis of precision metric using machine learning classifiers

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Multinomial Naïve Bayes	73.3	85.8	88.4	92.3	95.8	84.4	99.1	97.3	89.3
Passive Aggressive	94.5	90.0	87.6	98.7	96.9	97.1	96.6	93.6	91.1
Stochastic Gradient Descent	94.8	89.4	88.1	98.7	97.8	95.8	96.4	94.2	90.0
Logistic Regression	95.4	94.0	89.8	95.4	97.9	94.2	94.7	95.2	89.1
Support Vector Classifier	100	96.6	100	100	97.2	100	100	53.1	82.0
NuSVC	94.5	96.3	90.4	97.5	97.9	91.6	41.9	77.8	85.7
LinearSVC	96.1	88.8	87.8	98.1	96.5	96.2	93.6	94.1	89.4
Multi-Layer Perceptron	93.6	93.1	84.0	97.2	96.5	96.2	92.1	94.3	89.9
Decision Tree	80.2	81.9	75.0	97.2	96.6	92.8	87.1	88.3	80.1
AdaBoost	89.9	88.7	83.6	97.2	96.6	94.1	92	90.6	84.9
Gradient Boosting	91.9	92.6	89.0	98.4	97.8	95.4	90.7	90.4	89.9
Voting Classifier	95.8	94.2	89.2	98.8	97.8	96.8	96.8	95.0	89.9

Bold indicates the best performance metric values among the other models

Table 8 Comparative analysis of recall metric using machine learning classifiers

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Multinomial Naïve Bayes	95.9	91.5	79.7	94.6	95.4	94	77.9	85.4	85.0
Passive Aggressive	92.2	88.1	85.0	98.0	98.3	94.8	96.0	95.3	86.7
Stochastic Gradient Descent	91.7	90.9	84.9	97.5	96.7	95.6	96.2	95.7	89.7
Logistic Regression	87.7	88	84.2	97.7	98.5	93.7	95.5	96.5	90.7
Support Vector Classifier	48.2	65.7	48.2	52.9	67.4	52.9	52.5	96.3	82.3
NuSVC	72.4	75.9	82.6	88.7	80.2	94.2	99.4	95.8	91.4
LinearSVC	91.1	86.4	86	97.8	98.9	94.9	96.6	95.1	90.1
Multi-Layer Perceptron	91.9	89.9	84.1	97.2	98.2	93.8	96.4	95.8	89.8
Decision Tree	80.7	78.3	73.9	94.5	96.1	90.9	89.2	89.6	82.3
AdaBoost	83.5	81.8	80.9	97.5	96.8	92.6	93.6	94.5	87.8
Gradient Boosting	86.4	85.0	82.0	98.7	98.4	96.3	94.2	95.1	89.9
Voting Classifier	91.6	89.7	85.1	98.0	98.9	95.4	95.0	96.6	89.8

Bold indicates the best performance metric values among the other models

Table 9 Comparative analysis of specificity metric using machine learning classifiers

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Multinomial Naïve Bayes	79.6	87.5	88	91.6	95.2	84.3	98.7	96.4	87.4
Passive Aggressive	94.8	90.5	88.2	98.5	96.6	96.6	96.2	93.1	89.6
Stochastic Gradient Descent	95	90.4	88.6	98.5	97.5	95.3	96	93.7	88.9
Logistic Regression	95.3	94	89.9	94.9	97.4	93.4	94.2	94.8	88.2
Support Vector Classifier	0	94.4	0	0	95.1	0	0	65.3	80.2
NuSVC	92.9	95.4	90.2	96.9	96.5	90.9	61.3	79.6	85.2
LinearSVC	96.2	89.3	88.5	97.8	96.2	95.7	95.9	93.5	88.4
Multi-Layer Perceptron	93.9	93.4	85.1	96.9	96.1	95.6	94.3	93.2	86.1
Decision Tree	81.7	82.4	76.4	96.8	96.2	91.7	86.1	87.3	78.6
AdaBoost	89.9	88.6	84.2	96.9	96.2	93.2	91.3	90.1	83.9
Gradient Boosting	92	92.5	88.9	98.2	97.5	94.8	90.1	89.9	85.2
Voting Classifier	95.9	94.3	89.5	98.6	97.6	96.3	96.4	94.6	88.8

Bold indicates the best performance metric values among the other models

Table 10 Comparative analysis of *F1*-score metric using machine learning classifiers

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Multinomial Naïve Bayes	86.9	89.4	83.6	93.4	95.5	88.9	87.2	90.9	87.0
Passive Aggressive	93.4	89.2	86.5	98.3	97.5	95.9	96.2	94.4	88.8
Stochastic Gradient Descent	93.3	90.6	86.7	98.0	97.2	95.6	96.2	94.9	89.8
Logistic Regression	91.3	90.9	86.8	96.5	98.1	93.9	95.0	95.8	89.8
Support Vector Classifier	65	77.4	65	69.1	79.6	69.1	68.8	68.4	82.1
NuSVC	81.3	84.5	86.2	92.8	88.1	92.8	58.9	85.8	88.4
LinearSVC	93.5	87.8	87.2	97.9	97.3	94.9	95	94.5	89.7
Multi-Layer Perceptron	92.8	91.6	84.5	97.2	97.3	94.9	93.2	94.1	88.7
Decision Tree	81.1	80.2	75.1	95.8	96.3	91.8	88.1	88.9	81.1
AdaBoost	86.5	85.0	82.5	97.3	96.6	93.3	92.7	92.5	86.3
Gradient Boosting	89.1	88.5	85.3	98.5	98.0	95.8	92.4	92.6	89.9
Voting Classifier	93.7	91.8	87.1	98.3	98.3	96	95.8	95.7	89.8

Bold indicates the best performance metric values among the other models

The objective of ROC is to notice the increase in false positive rates (FPR) with an increase in true positive rates (TPR) with a varying threshold of the classifiers used in this paper. The performance of class models at various thresholds is shown through graphs in Fig. 12. The curve drawn in the graph is known as receiver operating characteristic (ROC) curve. The ROC curves for News trends, Kaggle and Reuters are plotted using two parameters such as FPR and TPR as given by Eqs. (15) and (16).

$$\text{TPR} = \frac{a}{a + c} \quad (15)$$

$$\text{FPR} = \frac{b}{b + d} \quad (16)$$

Here, a , b , c , and d represent the TP, FP, FN and TN rates, respectively. predicted document as 'Real' when actually it

is 'Fake,' known as false positive (FP) rate, c counts the predicted document as 'Fake' when actually it is 'Real,' known as false negative (FN) rate, d counts the predicted document as 'Real' when actually it is 'Real,' known as true negative (TN) rate.

6.1 Major findings

The major findings deal with the question of whether a trained classifier using old news articles can give accurate and efficient results to categorize the differences between fake and real contents. It has been observed that the performance of classifying news article depends on the corpus and type of classification model. In our experiment, three corpora have been collected from three different sources (<https://bit.ly/2zVRLxK>, <https://bit.ly/2Ex5VsX>, <https://bit.ly/2BmqBQE>). Each corpus is divided into training (0.67)

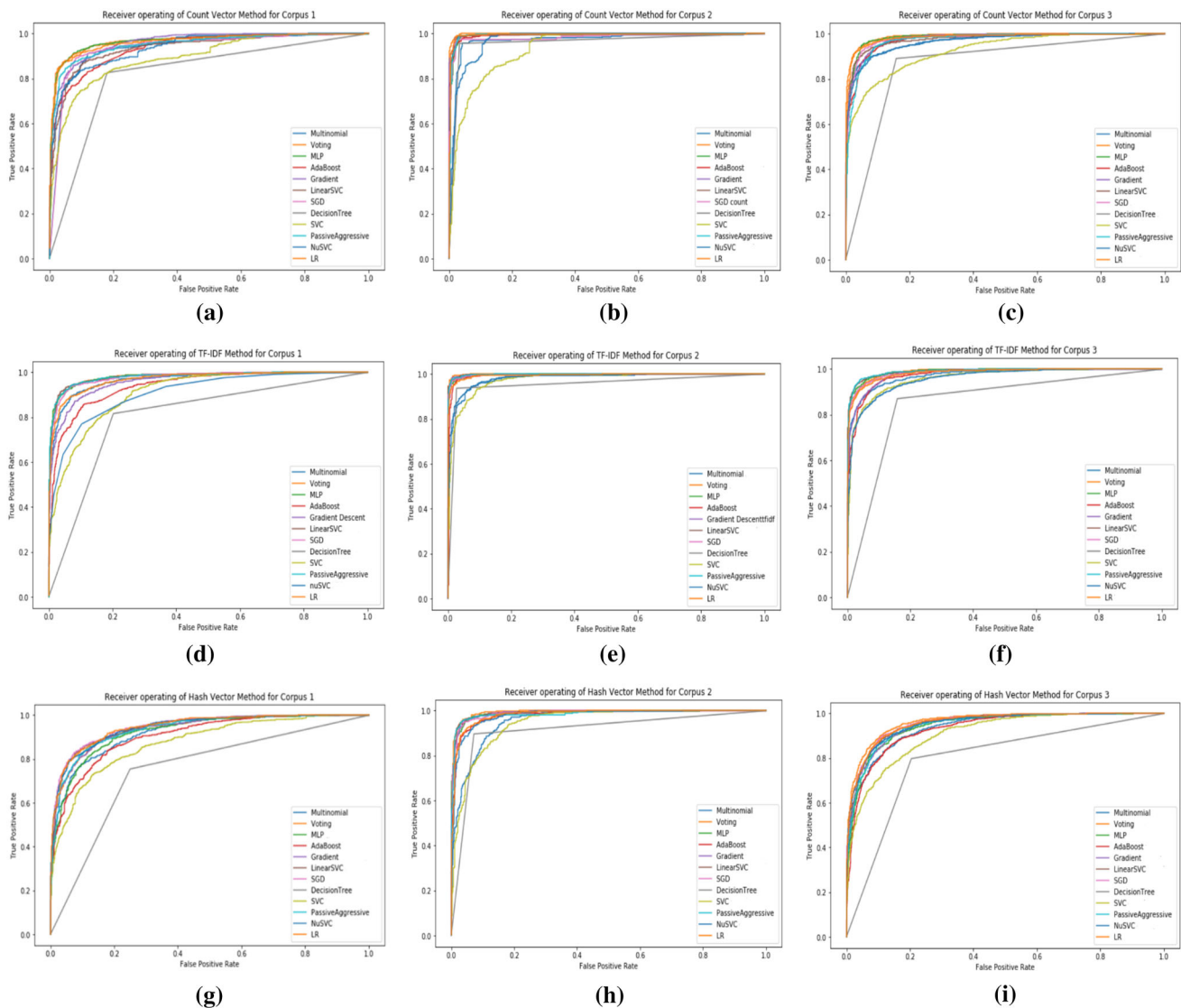


Fig. 12 TPR versus FPR at different classification thresholds for **a** Count-Vectorizer on News Trends, **b** Count-Vectorizer on Kaggle, **c** Count-Vectorizer on Reuters, **d** TF-IDF on News Trends, **e** TF-IDF on

Kaggle, **f** TF-IDF on Reuters, **g** Hashing-Vectorizer on News Trends, **h** Hashing-Vectorizer on Kaggle and **i** Hashing-Vectorizer on Reuters datasets

and testing (0.33) sets. The experiment was performed on these chosen datasets using Term Frequency-Inverse Document Frequency (TF-IDF), Count-Vectorizer (CV) and Hashing-Vectorizer (HV) feature extraction techniques. From the accuracy perspective, Passive Aggressive (93.2%) and LinearSVC (93.2%) outperform other models for all three (News Trends, Kaggle, Reuters) corpora, whereas the Passive Aggressive (96%) and the LinearSVC (95.9%) perform best using TF-IDF, Logistic Regression (94.9%) and Stochastic Gradient Descent (94.2%) perform best using CV, and LinearSVC (90.6%) and Stochastic Gradient Descent (90.5%) perform best using HV individually for all three corpora.

A classifier is considered usable only if it achieves both high precision and recall. To average out the results of both precision and recall, $F1$ -score is taken into consider-

ation. On evaluating $F1$ -score metric, it was observed that Passive Aggressive (93.3%), Stochastic Gradient Descent (93.5%) and LinearSVC (93%) outperform other models on all three news article corpora using TF-IDF, CV and HV feature extraction techniques. The Passive Aggressive (95.9%), Stochastic Gradient Descent (95.8%) and LinearSVC (95.4%) perform best using TF-IDF, Logistic Regression (94.9%), Stochastic Gradient Descent (94.2%), Passive Aggressive (93.7%), LinearSVC (93.2) perform best using CV, and Passive Aggressive (90.4%), Stochastic Gradient Descent (90.7%) and LinearSVC (90.6%) perform best using HV. The proposed multi-level voting model outperforms the Passive Aggressive model by 0.8%, 0.6%, and 1.0% using TF-IDF approach; outperforms the Logistic Regression by 2.6%, 0.7%, 0.8% using CV approach;

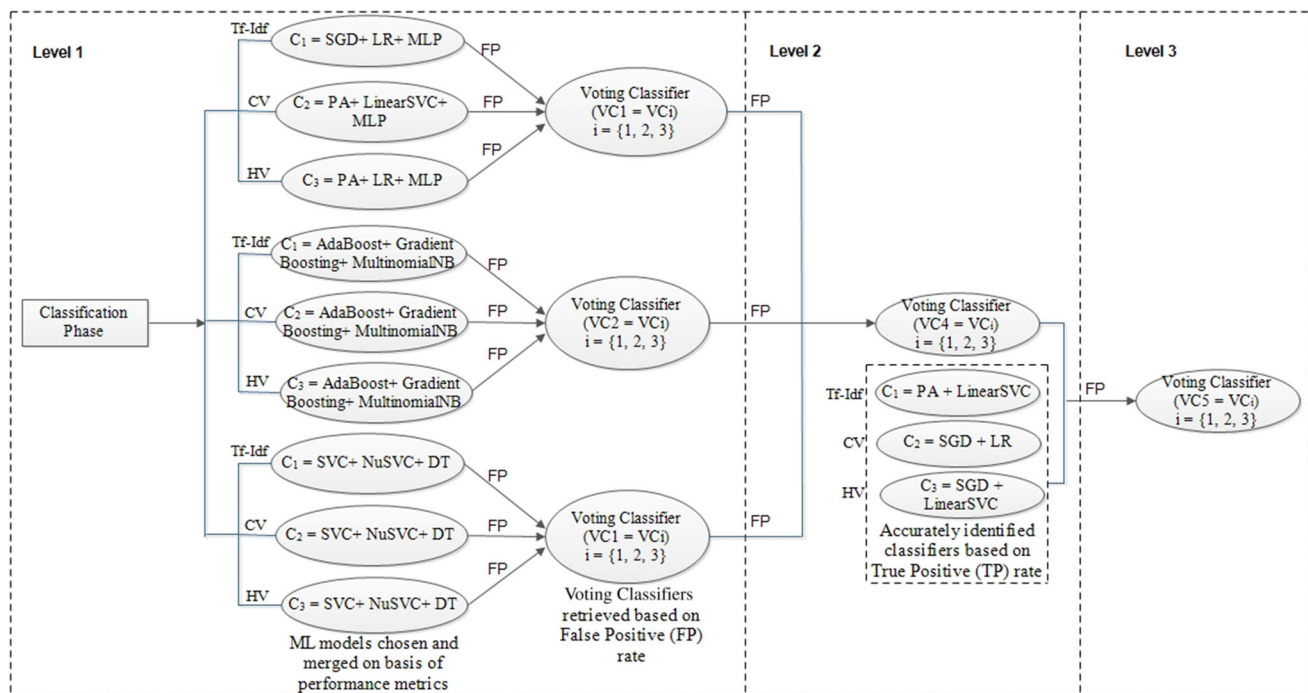


Fig. 13 Architecture of the proposed multi-level voting model

outperforms the LinearSVC by 0.0%, 0.5%, 0.9% using HV approach on News Trends, Kaggle and Reuters corpus, respectively.

To evaluate the predictive performance of our approach, ROC_AUC is plotted. Passive Aggressive, Stochastic Gradient Descent and LinearSVC, Gradient Boosting, Logistic Regression outperform other ML models on the basis of ROC_AUC metric with $\text{TPR} > 0.97$. Passive Aggressive, Stochastic Gradient Descent and LinearSVC gives $\text{TPR} > 0.97$ using TF-IDF, Logistic Regression gives $\text{TPR} > 0.98$ using CV, and Stochastic Gradient Descent, LinearSVC, Logistic Regression gives $\text{TPR} > 0.95$ for News Trends, Kaggle and Reuters datasets, respectively. Based on the training time required by various ML classifiers, it has been observed that there is a trade-off between efficiency and accuracy. The training time required by HV is less than TF-IDF and CV technique, but it compromises accuracy metric. It has been analyzed that the hashing technique is useful when the focus is to achieve high efficiency on a huge dataset. The two ML classifiers such as Logistic Regression and LinearSVC are chosen among other models which result in both high accuracy and efficiency to overcome the trade-off issue.

7 Proposed multi-level voting model

The proposed multi-level voting model not only helps to improve the accuracy but also helps to reduce the training time of the classifiers. The reduction in training time helps

to increase the efficiency of our model by introducing parallel methods where the base learners are generated parallelly. The motivation behind the proposed model is to analyze the independence between the base learners. Three levels are proposed to perform the experiment as discussed below.

Level 1 Sets of three ML classifiers are merged based on their performance metric (FP rate) to apply Voting Classifier. The voting models (VC1, VC2, VC3) are retrieved.

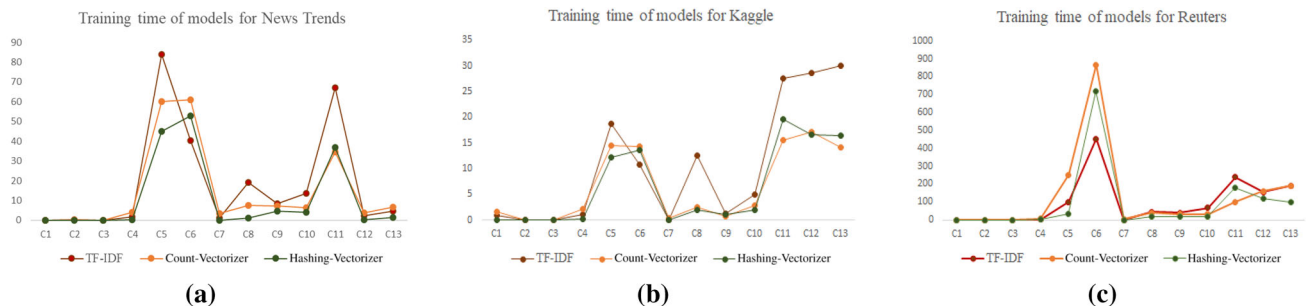
Level 2 A Voting Classifier (VC4) is retrieved after merging the three models (VC1, VC2, VC3) on the basis of their false positive rate.

Level 3 The false predictions from Voting Classifier (VC4) are merged with PA and LinearSVC for TF-IDF, LR and SGD for CV and SGD, LinearSVC for HV to get the final prediction.

On the basis of the minimum false positive (FP) rate, the ML models are merged to overcome the weakness of the existing individual models. The minimum the FP ratio, the more accurate the model will be to predict the content as fake. Three feature extraction techniques (TF-IDF, CV and HV) are used to extract the features from a collected dataset. Based on the FP ratio, the models are selected and merged to give an appropriate prediction. First cluster at level 1, SGD (News Trends), LR (Kaggle) and MLP (Reuters) using TF-IDF; SGD (News Trends), LinearSVC (Kaggle) and MLP (Reuters) using CV; PA (News Trends), LR (Kaggle) and MLP (Reuters) using HV are merged to build the Voting Classifier (VC1). Second cluster at level 1, AdaBoost, Gradient Boosting and MultinomialNB (News Trends, Kaggle

Table 11 Comparative analysis of the proposed multi-level voting model with Voting Classifier

Models	News Trends			Kaggle			Reuters		
	TF-IDF	CV	HV	TF-IDF	CV	HV	TF-IDF	CV	HV
Accuracy									
Voting classifier	93.8	92.1	87.3	98.3	98.3	95.9	96.1	96.4	90.3
Multi-Voting Classifier	94.3	93.6	87.1	98.9	98.7	95.8	97.2	96.5	90.2
Precision									
Voting Classifier	95.8	94.2	89.2	98.8	97.8	96.8	96.8	95	89.9
Multi-Voting Classifier	96.4	94.3	89.6	99.1	98.3	96.8	98.4	97.6	90.4
Recall									
Voting Classifier	91.6	89.7	85.1	98	98.9	95.4	95	96.6	89.8
Multi-Voting Classifier	93.1	91.4	85.2	98.7	98.8	94.4	96.8	95.2	90.8
Specificity									
Voting Classifier	95.9	94.3	89.5	98.6	97.6	96.3	96.4	94.6	88.8
Multi-Voting Classifier	96	92.4	86.1	98.2	97.1	93.6	95.7	94.6	87.9
F1-score									
Voting Classifier	93.7	91.8	87.1	98.3	98.3	96	95.8	95.7	89.8
Multi-Voting Classifier	94.7	92.8	87.3	98.8	98.7	95.5	97.7	97	90.5

**Fig. 14** Training time comparison for **a** News Trends, **b** Kaggle and **c** Reuters dataset on the basis of TF-IDF, Hashing-Vectorizer and Count-Vectorizer feature extraction techniques

and Reuters) using TF-IDF, CV and HV are merged to build VC2. Third cluster at level 1, SVC, NuSVC and DT (News Trends, Kaggle and Reuters) using TF-IDF, CV and HV are merged to build VC3 as shown in Fig. 13.

Based on FP ratio of three proposed Voting Classifier (VC1, VC2 and VC3), a fourth Voting Classifier (VC4) is retrieved at level 2 based on the FP rates of VC1, VC2 and VC3 classifiers. At level 3, PA, LinearSVC using TF-IDF; SGD and LR using CV; SGD and LinearSVC using HV are retrieved based on TP rate and are further clustered with VC4. Based on FP rate, VC5 is retrieved to give the final prediction of the proposed model.

8 Evaluating the performance of multi-level voting model

In proposed multi-level Voting Classifier, the best three ML models are combined from each feature extraction technique

as discussed in the previous section. It can be observed from Table 11 that the proposed model outperforms Voting Classifier by 0.73%, 0.66% and 0.13% using TF-IDF, CV and HV, respectively, in terms of accuracy metric. Similarly, the proposed model also gives significant improvement for precision, recall, specificity and F1-score performance measures.

To compare the vectorization methods used in terms of training time required to train the data, Fig. 14 depicts that the HV technique is more efficient than other feature extraction methods as the time required to train the proposed multi-level voting model is least for News Trends and Reuters datasets. To achieve better efficiency (less training time), TF-IDF method is only suitable when data are not too huge, whereas hashing can be used for huge datasets when there is a requirement to make a trade-off between efficiency and accuracy. After performing our experiment, some conclusions are drawn in next section which helps to analyze the best classifier to be chosen for both high efficiency and accuracy.

Table 12 Comparison of the existing models with a proposed system for Reuters corpus

Authors	Proposed approach	Features	Model accuracy
Stein and Zu Eissen (2008)	Presented a Bayesian classification approach using class-specific features for automatic text classification	Topic-based and document-based modeling	LapPLSI—74.6%
Mishu and Rafiuddin (2016)	Classified text document using various classifiers	Document frequency	MultinomialNB—72% LR—73.5% SGD—76% SVC—78% LinearSVC—83.3% Voting—89% SGD—96.2%
Analysis based on individual ML models	Classified articles as fake or real using various classifiers	TF-IDF	PA—96.2% LinearSVC—96.3% MLP—96.2% AdaBoost—92.5% Voting—96.4% MLP—96.2% DT—88.5% LR—95.7% MultinomialNB—89.9% Gradient Boosting—92.5% SVC—81.3% NuSVC—88.3%
Proposed multi-level voting model	Fake news detection system is proposed to achieve high accuracy and high efficiency	Count-Vectorizer	Multi-level voting—97.2%
		Hashing-Vectorizer	Multi-level voting—96.5%
		Hashing-Vectorizer	Multi-level voting—90.2%

9 Comparison with the existing works

The results given by our system are also compared with other existing works on fake news detection as shown in Tables 12, 13 and 14 using Reuters, Kaggle and News Trends corpora, respectively. Researchers have used ML approaches to perform fake news detection on various social media platforms. It has been analyzed from Table 12 that the feature extraction technique based on document frequency used by Mishu and Rafiuddin (2016) for MultinomialNB (72%), SVC (78%) and Voting Classifier (89%) (Mishu and Rafiuddin 2016) does not give better accuracy than our proposed system for all three classifiers. In our proposed system, MultinomialNB (89.9%), SVC (81.3%) and Voting (96.4%) outperform Mishu and Rafiuddin (2016) models.

It has been observed from Table 13 that the feature extraction technique based on TF-IDF used by Ahmed et al. (2017) for ML classifiers does not give better accuracy than our proposed system. The recorded observation shows that our

proposed multi-level voting model (97.2%) outperforms Gradient Boosting by 0.4% using TF-IDF, LR by 2.3% using CV and NuSVC by 3.6% using HV feature extraction techniques, respectively, for Kaggle corpus. From Table 14, it can be analyzed that our proposed multi-level voting model outperforms PA by 0.8% using TF-IDF, MultinomialNB by 4.3% using CV and NuSVC by 0.8% using HV feature extraction techniques, respectively, for News Trends corpus. From the comparative analysis, it has been analyzed that our proposed multi-level voting model using all three feature extraction technique outperforms when compared to the individual performance metrics of ML models used by Ahmed et al. (2017) and Mishu and Rafiuddin (2016) for developing a system for automatic classification of news article features to label them as fake or real news article. It has been observed from Table 12 that the proposed model outplays the PA model by 0.9% using TF-IDF, LR model by 0.8% using CV and NuSVC model by 1.9% using HV feature extraction techniques, respectively, when compared with their performance metrics.

Table 13 Comparison of the existing models with a proposed system for Kaggle corpus

Authors	Proposed approach	Features	Model accuracy
Ahmed et al. (2017)	Proposed a fake news detection model using n-gram analysis and ML techniques	n-gram based TF and TF-IDF	LSVM—92% KNN—83.1% SVM—86% DT—89% SGD—89% LR—89%
Analysis based on individual ML models	Classified articles as fake or real using various classifiers	Count-Vectorizer	MultinomialNB—93.2% SVC—52.9% LR—96.4% MLP—97.1% DT—95.6% Voting—98.3%
		TF-IDF	LinearSVC—97.9% SGD—98% PA—98.3% AdaBoost—97.3% Gradient Boosting—98.5% Voting—98.3%
Proposed multi-level voting model	Fake news detection system is proposed to achieve high accuracy and high efficiency	Hashing-Vectorizer TF-IDF	NuSVC—92.2% Multi-level voting—98.9%
		Count-Vectorizer	Multi-level voting—98.7%
		Hashing-Vectorizer	Multi-level voting—95.8%

10 Conclusion and future scope

The goal in this paper is to analyze the best known supervised technique for detecting fake news. When dealing with the machine learning classifiers, the key question is not to find a learning classifier superior to others but rather to find the conditions under which particular model outperform others for a given problem. The set of attributes extracted from the corpus taken uses three feature extraction techniques (TF-IDF, CV and HV) to feed the extracted feature vectors into the selected machine learning models. Some characteristics taken from datasets for learning task are categorical attributes, missing values, headlines of the article, the body of article and the publisher name. Various classifiers, Multinomial Naïve Bayes (MultinomialNB), Passive Aggressive (PA), Stochastic Gradient Descent (SGD), Logistic Regression (LR), Support Vector classifier (SVC), NuSVC, LinearSVC, Multi-Layer Perceptron (MLP), Deci-

sion Tree (DT), AdaBoost, Gradient Boosting and Voting Classifier were analyzed on the basis of performance measures. After analyzing the classifiers, the focus was to utilize the strengths of one model to complement the weakness of another. So, the multi-level voting model was proposed, which integrates various ML models based on their FP rates to retrieve a news Voting Classifier to retrieve better prediction analysis. The developed model helps to solve the trade-off issue between accuracy and efficiency.

In the future, a web-based GUI will be created for the proposed fake news detection system to classify the news as fake or real on real-time social media platforms such as Facebook, Instagram, Twitter and WhatsApp. Also, the annotated dataset in form of images (with textual content written on them) will be collected and maintained from Facebook and Reddit platforms. The annotated dataset can be used for detecting fake images in the future as no such dataset is available at present. The proposed system has the potential

Table 14 Comparison of the existing models with a proposed system for News Trends corpus

Authors	Proposed approach	Features	Model accuracy
Kuleshov et al. (2018)	Author shows the existence of adversarial examples in natural language classification	n-gram	NB—93%
Analysis based on individual ML models	Classified articles as fake or real using various classifiers	TF-IDF	PA—93.5%
			SGD—93.4%
			LinearSVC—93.6%
			MLP—93%
			DT—81.3%
			AdaBoost—86.7%
			Gradient Boosting—89.2%
			Voting—93.8%
			LR—91.4%
			Count-Vectorizer
Proposed multi-level voting model	Fake news detection system is proposed to achieve high accuracy and high efficiency	TF-IDF	MultinomialNB—89.3%
			SVC—74.1%
			Hashing-Vectorizer
			NuSVC—86.3%
			Multi-level voting—94.3%
		Count-Vectorizer	Multi-level voting—93.6%
			Hashing-Vectorizer
			Multi-level voting—87.1%

to provide an impulse to various emerging applications such as controlling the spread of fake news during elections, terrorism, natural calamities, crimes for the betterment of the society.

Acknowledgements This Publication is an outcome of the R&D work undertaken in the project under the Visvesvaraya PhD Scheme of Ministry of Electronics and Information Technology, Government of India, being implemented by Digital India Corporation (formerly Media Lab Asia).

Funding Funding was provided by Digital India Corporation (formerly Media Lab Asia) (Grant No. U72900MH2001NPL133410).

Compliance with ethical standards

Conflict of interest The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- Aggarwal A, Rajadesingan A, Kumaraguru P (2012) PhishAri: automatic realtime phishing detection on twitter. In: eCrime researchers summit (eCrime). IEEE, pp 1–12
- Aggarwal A, Kumar S, Bhargava K, Kumaraguru P (2018) The follower count fallacy: detecting twitter users with manipulated follower count
- Ahmed F, Abulaish M (2012) An MCL-based approach for spam profile detection in online social networks. In: IEEE 11th international conference on trust, security and privacy in computing and communications (TrustCom). IEEE, pp 602–608
- Ahmed H, Traore I, Saad S (2017) Detection of online fake news using n-gram analysis and machine learning techniques. In: International conference on intelligent, secure, and dependable systems in distributed and cloud environments. Springer, pp 127–138
- Alahmadi A, Joorabchi A, Mahdi AE (2013) A new text representation scheme combining bag-of-words and bag-of-concepts approaches for automatic text classification. In: 2013 7th IEEE GCC conference and exhibition (GCC). IEEE, pp 108–113
- Batchelor O (2017) Getting out the truth: the role of libraries in the fight against fake news. *Ref Serv Rev* 45(2):143
- Benevenuto F, Rodrigues T, Almeida V, Almeida J, Gonçalves M (2009) Detecting spammers and content promoters in online video social networks. In: Proceedings of the 32nd international ACM SIGIR conference on research and development in information retrieval. ACM, pp 620–627
- Benevenuto F, Magno G, Rodrigues T, Almeida V (2010) Detecting spammers on twitter. In: Collaboration, electronic messaging, anti-abuse and spam conference (CEAS), vol 6, p 12
- Caetano JA, de Oliveira JF, Lima HS, Marques-Neto HT, Magno G, Meira W Jr, Almeida VA (2018) Analyzing and characterizing political discussions in WhatsApp public groups. *arXiv preprint arXiv:1804.00397*
- Canini KR, Suh B, Pirolli PL (2011) Finding credible information sources in social networks based on content and social structure. In: IEEE third international conference on social computing (SocialCom). IEEE third international conference on privacy, security, risk and trust (PASSAT). IEEE, pp 1–8

- Chen Y, Conroy NJ, Rubin VL (2015) Misleading online content: recognizing clickbait as false news. In: Proceedings of the 2015 ACM on workshop on multimodal deception detection. ACM, pp 15–19
- Chhabra S, Aggarwal A, Benevenuto F, Kumaraguru P (2011) Phi.sh\$ocial: the phishing landscape through short URLs. In: Proceedings of the 8th annual collaboration, electronic messaging, anti-abuse and spam conference. ACM, pp 92–101
- Conroy NJ, Rubin VL, Chen Y (2015) Automatic deception detection: methods for finding fake news. *Proc Assoc Inf Sci Technol* 52(1):1
- D'Angelo G, Palmieri F, Rampone S (2019) Detecting unfair recommendations in trust-based pervasive environments. *Inf Sci* 486:31
- Dewan P, Kumaraguru P (2015) Towards automatic real time identification of malicious posts on facebook. In: 13th Annual conference on privacy, security and trust (PST). IEEE, pp 85–92
- Dewan P, Kumaraguru P (2017) Facebook inspector (FbI): towards automatic real-time detection of malicious content on Facebook. *Soc Netw Anal Min* 7(1):15
- Dewan P, Gupta M, Goyal K, Kumaraguru P (2013) Multiosn: realtime monitoring of real world events on multiple online social media. In: Proceedings of the 5th IBM collaborative academia research exchange workshop. ACM, p 6
- Fake news on whatsapp. <http://bit.ly/2miuv9j>. Last accessed 27 Aug 2019
- Gao H, Hu J, Wilson C, Li Z, Chen Y, Zhao BY (2010) Detecting and characterizing social spam campaigns. In: Proceedings of the 10th ACM SIGCOMM conference on internet measurement. ACM, pp 35–47
- Garimella K, Tyson G (2018) WhatsApp, doc? A first look at WhatsApp public group data. arXiv preprint [arXiv:1804.01473](https://arxiv.org/abs/1804.01473)
- Gupta A, Kumaraguru P (2012a) Credibility ranking of tweets during high impact events. In: Proceedings of the 1st workshop on privacy and security in online social media. ACM, p 2
- Gupta A, Kumaraguru P (2012b) Twitter explodes with activity in Mumbai blasts! a lifeline or an unmonitored daemon in the lurking? Technical report
- Gupta A, Lamba H, Kumaraguru P (2013a) \$ 1.00 per rt #Boston-Marathon #PrayForBoston: analyzing fake content on twitter. In: eCrime researchers summit (eCRS). IEEE, pp 1–12
- Gupta A, Lamba H, Kumaraguru P, Joshi A (2013b) Faking sandy: characterizing and identifying fake images on twitter during hurricane sandy. In: Proceedings of the 22nd international conference on world wide web. ACM, pp 729–736
- Jain P, Kumaraguru P (2016) On the dynamics of username changing behavior on twitter. In: Proceedings of the 3rd IKDD conference on data science. ACM, p 6
- Kaggle database. <https://bit.ly/2BmqBQE>. Last accessed 22 Oct 2017
- Kaggle database. <https://bit.ly/2Ex5VxX>. Last accessed 24 Oct 2017
- LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521(7553):436
- Kuleshov V, Thakoor S, Lau T, Ermon S (2018) Adversarial examples for natural language classification problems
- Magdy A, Wanas N (2010) Web-based statistical fact checking of textual documents. In: Proceedings of the 2nd international workshop on search and mining user-generated contents. ACM, pp 103–110
- Markines B, Cattuto C, Menczer F (2009) Social spam detection. In: Proceedings of the 5th international workshop on adversarial information retrieval on the web. ACM, pp 41–48
- Mishu SZ, Rafiuddin S (2016) Performance analysis of supervised machine learning algorithms for text classification. In: 19th International conference on computer and information technology (ICCIT). IEEE, pp 409–413
- News trends database. <https://bit.ly/2zVRLxK>. Last accessed 18 Oct 2017
- Pontes T, Magno T, Vasconcelos M, Gupta A, Almeida J, Kumaraguru P, Almeida V (2012a) Beware of what you share: inferring home location in social networks. In: IEEE 12th international conference on data mining workshops (ICDMW). IEEE, pp 571–578
- Pontes T, Vasconcelos M, Almeida J, Kumaraguru P, Almeida V (2012b) We know where you live: privacy characterization of foursquare behavior. In: Proceedings of the 2012 ACM conference on ubiquitous computing. ACM, pp 898–905
- Qazvinian V, Rosengren E, Radev DR, Mei Q (2011) Rumor has it: identifying misinformation in microblogs. In: Proceedings of the conference on empirical methods in natural language processing. Association for Computational Linguistics, pp 1589–1599
- Rubin VL, Conroy NJ, Chen Y (2015) Towards news verification: deception detection methods for news discourse. In: Hawaii international conference on system sciences
- Rubin V, Conroy N, Chen Y, Cornwell S (2016) Fake news or truth? Using satirical cues to detect potentially misleading news. In: Proceedings of the second workshop on computational approaches to deception detection, pp 7–17
- Ruchansky N, Seo S, Liu Y (2017) CSI: a hybrid deep model for fake news detection. In: Proceedings of the 2017 ACM on conference on information and knowledge management. ACM, pp 797–806
- Sen I, Aggarwal A, Mian S, Singh S, Kumaraguru P, Datta A (2018) Worth its weight in likes: towards detecting fake likes on Instagram. In: Proceedings of the 10th ACM conference on web science. ACM, pp 205–209
- Shu K, Sliva A, Wang S, Tang J, Liu H (2017) Fake news detection on social media: a data mining perspective. *ACM SIGKDD Explor Newsl* 19(1):22
- Sirajudeen SM, Azmi NFA, Abubakar AI (2017) Online fake news detection algorithm. *J Theor Appl Inf Technol* 95(17):4114
- Stein B, Zu Eissen SM (2008) Retrieval models for genre classification. *Scand J Inf Syst* 20(1):3
- Volkova S, Shaffer K, Jang JY, Hodas N (2017) Separating facts from fiction: linguistic models to classify suspicious and trusted news posts on twitter. In: Proceedings of the 55th annual meeting of the association for computational linguistics (volume 2, short papers), vol 2, pp 647–653
- Wang AH (2010) Don't follow me: spam detection in twitter. In: Proceedings of the 2010 international conference on security and cryptography (SECRYPT). IEEE, pp 1–10
- Wei W, Wan X (2017) Learning to identify ambiguous and misleading news headlines. arXiv preprint [arXiv:1705.06031](https://arxiv.org/abs/1705.06031)
- Weimer M, Gurevych I, Mühlhäuser M (2007) Automatically assessing the post quality in online discussions on software. In: Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions. Association for Computational Linguistics, pp 125–128