# Hybrid Text Classification Method for Fake News Detection

**Prabhjot Kaur, Rajdavinder Singh Boparai, and Dilbag Singh**

*Abstract:Fake news will be news, stories or scams made to purposely misguide or delude perusers. As a rule, these accounts are made to impact individuals' perspectives, push a political motivation or cause disarray and can regularly be a gainful business for online distributers. Fake news stories can swindle individuals by looking like believed sites or utilizing comparative names and web delivers to trustworthy news associations. The fake news detection has the three phases which are pre-processing, feature extraction and classification. In the previous time Support Vector Machine (SVM) classification is applied for the fake news detection. To improve accuracy of the fake news hybrid classification model is designed in this research work. The proposed model is implemented in Python and results are analyzed in terms of accuracy, precision and recall. Experimental analysis shows that the proposed method outperforms competitive techniques.*

*Keywords: Fake news, SVM, Hybrid classifier*

## I. INTRODUCTION

"False reports" are used to symbolize untrue rumor or misinformation encompassing propaganda conversed with the help of conventional medium such as publishing houses, television and also through non-conventional medium such as social networking sites. The universal reasons of spreading such kind of reports engage confusing the readers; damaging the status of some individual, or for the attainment of some aggrandizement [1]. Fake news is considered as furthermost intimidation to democratic system, liberated discussion and the Western sort. Counterfeit reports are capable of producing damaging effects on persons and the civilization. Because of the fake news individuals may be deceived and admit bogus attitude. Counterfeit reports can alter the people's reaction towards factual reports. The reliability of whole report environment can be damaged through extensive broadcasting of bogus information [2]. Thus, the recognition of bogus information on social networking sites is necessary. False reports are deliberately printed to misinform customers; therefore the recognition normal news becomes difficult. The exploration of supplementary data from dissimilar perceptions is normal and essential for the development of an effectual and sensible bogus reports recognition scheme.

False reports discovery has freshly fascinated a mounting attention from the common community as well as researchers because of the online increasing of propaganda movement predominantly in medium outlets like social networking sites, supplies, report blogs, and online correspondents [3]. Some techniques offer immense assurance for developer for the construction of schemes which can involuntarily perceive counterfeit reports. Though, recognition of false information is a demanding duty to achieve because of its requirement of mock-up for the summarization of reports. For the classification of false report, a comparison is performed between genuine and fraudulent information. Furthermore, the assignment of evaluation of proposed news with the genuine report is a intimidating chore because of its prejudice and estimation [4]. A dissimilar method for the detection of bogus information is from stance recognition .attitude recognition is the procedure of mechanically noticing the association amid two portions of content. A way is explored in this research for the prediction of the deportment in association with information editorial and information caption couple [5]. Relying on the similarity between same news editorial text and caption, the posture amid them can be depicted as "ready", 'not ready', 'discussion' or 'neutral'. A lot of experiments are conducted with some conventional machine learning approaches for the settlement of a baseline and then comparison is performed between these results to the state-of-the art deep networks for the classification of the attitude amid editorial corpse and caption [6]. A number of phases are involved in the recognition of false reports. The ground of false information recognition is a comparatively novel region of investigation [7]. Hence, a small number of unrestricted data samples are available. A primarily new data sample is collected by the researcher through compilation of publicly obtainable information editorials. Information preprocessing prepares unprocessed information for supplementary dispensation [8]. The conventional information preprocessing technique initiates with data which is implicitly prepared for investigation without any response and convey the method of information compilation. In tokenization, the stream of content is breached into language, idioms, cryptogram or additional significant rudiments described as tokens. The aspiration of this technique is the examination of the vocabulary in a statement. The stream of symbols is converted into input for supplementary dispensation like parsing or passage withdrawal. Stemming is the procedure of converting the alternative patterns of a statement into an ordinary depiction called stem. For example, the words: "presentation", "presented", "presenting" could all be summarized to an ordinary illustration called "present". This procedure is utilized globally in content dispensation for performing data repossession [9].

This technique is relied on the assumption that affectation of a question with the expression presentation entails a curiosity in papers comprising the presentable and accessed words. Feature extraction addresses the difficulty related to the discovery of mainly compressed and educational pattern of characteristics for improving the effectiveness of information storage and dispensation is addressed by feature extraction [10]. Representation of characteristic vectors remnants the most ordinary and suitable mean of information depiction for categorization and deterioration issues. The preceding stride in the categorization procedure is the training of classifier. Dissimilar classifiers are used for the prediction of the group of the papers [11]. After the training of classifier, four assessment outcomes like accuracy, remembrance, F-measure and exactness are obtained which are categorized in fake and genuine news. Fake news can be described in the form of a made-up story with a purpose of deceiving someone. Fake news is a experience which is having a noteworthy bang on our common life, particularly in the political world.

The remaining portion of this paper contains the "Literature Review" section that has the table of comparison. The part Research Methodology explores the evaluation of data pre-processing, feature extraction and classification. The next part contains result and discussion. The concluding portion is examined in the section "Conclusion".

## II.    LIERATURE REVIEW

Katherine Clayton, et al. (2019) evaluated the effectiveness of different technologies using which the false stories posted on various social networking sites like Facebook could be identified [12]. The belief in false news is reduced to a moderate level by the "Disputed" and "Rated false" tags. It has seen that "Disputed" approach provided higher accuracy results in previous approaches. However, this research showed that the belief in misinformation for reduced to greater extent by the "Rated False" mechanism. Simulation results achieved showed that the effective approaches derived from this work proved to be highly beneficial when applied in real time scenarios.

Atodiresei, et al. (2018) suggested a scheme for the identification of fraudulent twitter followers and fraudulent twitter reports [13]. The presented approach returned a pattern of figures about the authenticity of tweets. By the means of tweeter groups and tweeter content, the projected approach did not achieve its main objective yet. Recognizing certainly that the presented report was false or not relied exclusively on its reputation in the similar social media platform was not a good suggestion. For the identification of false reports, individual possessions were used by the face book for the investigation of reputed news.

Aldwairi, et al. (2018) defined an easy and efficient methodology for the clients which allowed them for the installation of a simple technique into their individual account [14]. The presented tool also allowed users to utilize this tool for identification and elimination of probable click enticements. The proposed approach showed tremendous performance in the recognition of false report origins.

Girgis, et al. (2018) stated that the main aim of this study was the development of a classifier for the prediction of fraudulent reports [15]. The fakeness of a report was based on the text. Thus for the encounter of false report issue, a completely deep neural approach with the outlook of RNN scheme representation and LSTMs was implemented in the proposed study. The tested results depicted that GRU showed best accuracy in comparison with several other approaches.

Al-Ash, et al. (2018) proposed a research work for the modeling of vectors for the accommodation of false reports features [16]. The vector modeling was performed before the supplementary progression by speech approaches with the help of the Indonesian communication. The major objective of the projected research was the identification of false reports. With the help of support vector machine approach, the frequency was conversed in the form of ten-fold cross corroboration. A much admirable performance had been shown by the vector demonstration which utilized the phrase frequency.

Vedova, et al. (2018) proposed a new ML false reports recognition methodology that combined the characteristics of public and reports texts [17]. The development of the projected approach was performed with the utilization of content relied and public relied methodology. For the validation of proposed approach, a number of experiments were carried out. The tested results verified the authenticity of the proposed approach. The combination of both of these approaches was relied on threshold imperative. The threshold rule was capable of capturing the distinct assistances of the proposed approaches and also performed better in comparison with several other methodologies. In future, the training of classifier with opinion reality in other speeches will be performed for the extension of proposed approach in several other nations.

### A.    Table of Comparison

| Author | Year | Work Done | Pros and Cons |
|---|---|---|---|
| Katherine Clayton | 2019 | The effectiveness of different technologies was evaluated using which the false stories posted on various social networking sites like Facebook could be identified. | Simulation results achieved showed that the effective approaches derived from this work proved to be highly beneficial when applied in real time scenarios. However, the effect of social endorsements or other contextual cues on belief in false news articles was not tested. |
| Atodiresei | 2018 | A scheme was proposed for the identification of fraudulent twitter followers and fraudulent twitter reports. The presented approach returned a pattern of figures about the authenticity of tweets. | For the identification of false reports, individual possessions were used by the face book for the investigation of reputed news. However, limited sources were considered to collect data in this research due to which the complexity of data was not very high. |
| Aldwairi | 2018 | An easy and efficient methodology was proposed for the clients who allowed them for the installation of a simple technique into their individual account. | The performed test results showed a precision rate of 99.4% with the help of logistic approach. However, its performance was not tested in the presence of new datasets. |

| Girgis | 2018 | The major objective of the projected research was the identification of false reports. With the help of support vector machine approach, the frequency was conversed in the form of ten-fold cross corroboration. | The precision figures obtained by other approaches were 0.2166 and 0.215. However, this research did not focus on improving the accuracy. |
|---|---|---|---|
| Al-Ash | 2018 | The major objective of the projected research was the identification of false reports. With the help of support vector machine approach, the frequency was conversed in the form of ten-fold cross corroboration. | This was capable for the detection of novel features properly with the accuracy rate of 96.74%. For this study, total 2516 papers were utilized during axiom recognition and named unit detection procedure. Solutions to cases where the content of the news is true but the title or the comment to the content is misleading or click bait were not considered here. |
| Vedova | 2018 | A new ML false reports recognition methodology was proposed that combined the characteristics of public and reports texts. | The threshold rule was capable of capturing the distinct assistances of the proposed approaches and also performed better in comparison with several other methodologies. |

## III. RESEARCH METHODOLOGY

This research work is related to fake news detection. The fake news detection process has the various phases which are described below:-

A. **Data Pre-Processing:-** In this phase, the dataset is taken as input from the kaggle. In the input dataset no missing value is there and the input dataset will be tokenized. The tokenized dataset will be processed and unwanted information will be removed from the dataset.

B. **Feature Extraction: -** In the second phase, the k fold method is applied for the feature extraction. In the feature extraction phase, the relationship will be established between the targets set.

C. **Classification: -** In this phase, the input dataset will be divided into training and testing. The training dataset will be 60 percent of the whole data and 40 percent will be test dataset.

The KNN is the k nearest neighbor algorithm which can calculate the nearest neighbor values in the input dataset. The nearest neighbor values of the dataset are calculated using the Euclidian distance formula. The number k value is selected from the network and based on k number of values the data can be classified into certain classes. The number of hyper planes depends upon the number of classes into which data needs to be classified. The random forest highly efficient algorithm which can provide greats results even without the presence of hyper-parameter tuning is called random forest algorithm. Due to its high simplicity and the fact that both classification and regression tasks can use it this classifier is gaining huge popularity. In the training time, multitude of decision tree is calculated and the mean prediction of individual trees is given as output.
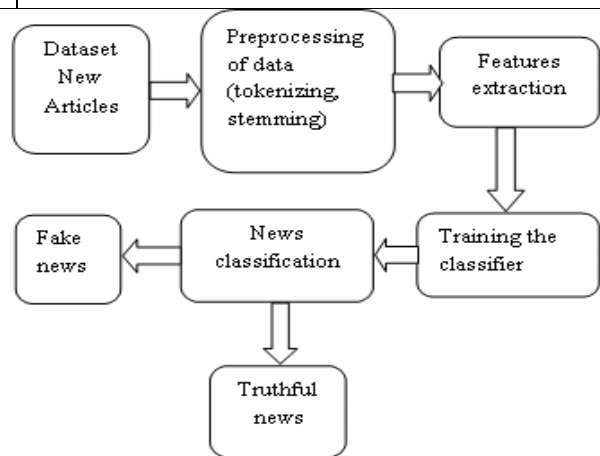


**Figure 1: Fake News Detection General Process**

## IV. RESULT AND DISCUSSION

The dataset of the fake news detection is collected from the kaggle. The dataset does not have any missing or redundant or missing values. The performance of the proposed model is tested in terms of certain parameters which are precision, recall, and accuracy.
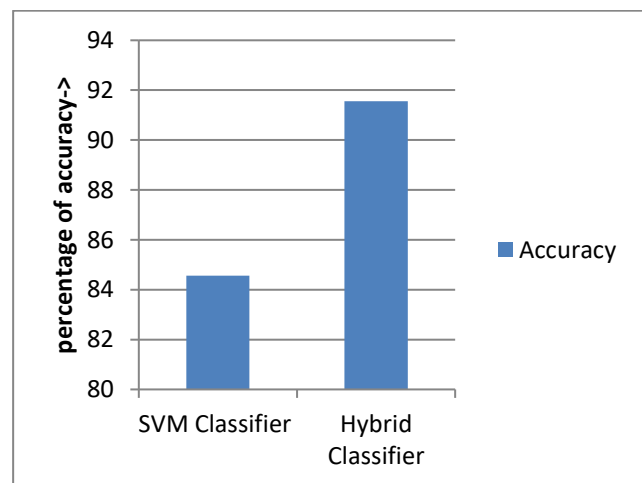


**Fig 2: Accuracy Analysis**

As shown in figure 2, the accuracy of the SVM classifier is compared with the hybrid classification. The hybrid classification model is the combination of KNN and random forest tree. It is analyzed that when the hybrid classification model is used accuracy is increased up to 8 percent.
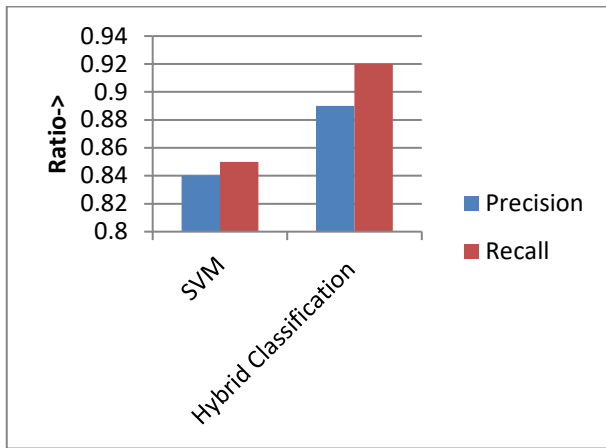


**Fig 3: Precision-Recall Analysis**

As shown in figure 3, the precision-recall value of the SVM classifier is compared with the hybrid classification. The precision-recall value of the hybrid classifier is high as compared to SVM classifier for the fake news detection.
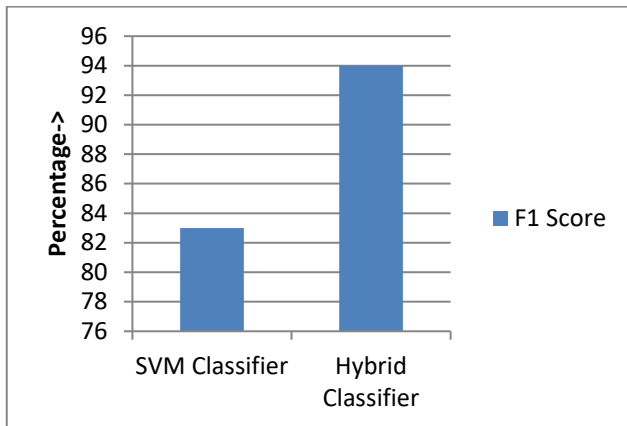


**Fig 3: F1 Score Analysis**

As shown in figure 3, the F1 score of SVM classifier is compared with hybrid classifier for the performance analysis. The hybrid classifier has high f1 score as compared to SVM classifier

## V.    CONCLUSION

In this work, it is found that the fake news detection is the application of the prediction analysis. The fake news detection process has the three phases which are pro-processing, feature extraction and classification. The hybrid classification model is designed in this research work for the fake news detection. The hybrid classification is the combination of the KNN and random forest. The performance of the proposed model is analyzed in terms of accuracy, precision and recall. The overall up to 8% results are improved with the use of hybrid model for the fake news detection.

## REFERENCES

1. Nir Kshetri, Jeffrey Voas, "The Economics of Fake News", IEEE, IT Professional, 2017, Volume: 19, Issue: 6, Pages: 8 - 12
2. Roger Musson, "Views: The frost report: fake news is nothing new", 2017, IEEE, Astronomy & Geophysics, Volume: 58, Issue: 3, Pages: 3.10 - 3.10
3. Hal Berghel, "Oh, What a Tangled Web: Russian Hacking, Fake News, and the 2016 US Presidential Election", IEEE, Computer, 2017, Volume: 50, Issue: 9, Pages: 87 - 91
4. Hal Berghel, "Alt-News and Post-Truths in the "Fake News" Era", IEEE, Computer, 2017, Volume: 50, Issue: 4, Pages: 110 - 114
5. Hal Berghel, "Lies, Damn Lies, and Fake News", IEEE, Computer, 2017, Volume: 50, Issue: 2, Pages: 80 - 85
6. Sneha Singhania, Nigel Fernandez, "3HAN: A Deep Neural Network for Fake News Detection", 2017, Neural Information Processing: 24th International Conference, ICONIP 2017, Guangzhou, China
7. Shashank Gupta, Raghuveer Thirukovalluru, Manjira Sinha, Sandya Mannarswamy, "CIMTDetect: A Community Infused Matrix-Tensor Coupled Factorization Based Method for Fake News Detection", 2018, IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)
8. Stefan Helmstetter, Heiko Paulheim, "Weakly Supervised Learning for Fake News Detection on Twitter", 2018, IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)
9. Akshay Jain, Amey Kasbe, "Fake News Detection", 2018, IEEE International Students' Conference on Electrical, Electronics and Computer Sciences
10. Chandra Mouli Madhav Kotteti, Xishuang Dong, Na Li, Lijun Qian, "Fake News Detection Enhancement with Data Imputation", 2018, IEEE 16th Int. Conf. on Dependable, Autonomic & Secure Comp., 16th Int. Conf. on Pervasive Intelligence & Comp., 4th Int. Conf. on Big Data Intelligence & Comp., and 3rd Cyber Sci. & Tech. Cong.
11. Shaban Shabani, Maria Sokhn, "Hybrid Machine-Crowd Approach for Fake News Detection", 2018 IEEE 4th International Conference on Collaboration and Internet Computing
12. Katherine Clayton, Spencer Blair, Jonathan A. Busam, Samuel Forstner, John Glance, Guy Green, Anna Kawata, Akhila Kovvuri, Jonathan Martin, Evan Morgan, Morgan Sandhu, Rachel Sang, Rachel Scholz Bright, Austin T. Welch, Andrew G. Wolf, Amanda Zhou, Brendan Nyhan, "Real Solutions for Fake News? Measuring the Efectiveness of General Warnings and Fact Check Tags in Reducing Belief in False Stories on Social Media", 2019, Springer Science, Business Media, LLC, part of Springer Nature
13. Costel-Sergiu Atodiresei, Alexandru Tănăselea, Adrian Iftene , "Identifying Fake News and Fake Users on Twitter", 2018, International Conference on Knowledge Based and Intelligent Information and Engineering Systems,Belgrade, Serbia
14. Monther Aldwairi, Ali Alwahedi, "Detecting Fake News in Social Media Networks", 2018, the 9th International Conference on Emerging Ubiquitous Systems and Pervasive Networks
15. Sherry Girgis, Eslam Amer, Mahmoud Gadallah, "Deep Learning Algorithms for Detecting Fake News in Online Text", 2018, 13th International Conference on Computer Engineering and Systems (ICCES)
16. Herley Shaori Al-Ash, Wahyu Catur Wibowo, "Fake News Identification Characteristics Using Named Entity Recognition and Phrase Detection", 2018, 10th International Conference on Information Technology and Electrical Engineering (ICITEE)
17. Marco L. Della Vedova, Eugenio Tacchini, Stefano Moret, Gabriele Ballarin, "Automatic Online Fake News Detection Combining Content and Social Signals", 2018, Procedding of the 22nd Conference of Fruct Assoction.
18. M.Vasuki, J. Arthi, K. Kayalvizhi," Decision Making Using Sentiment Analysisfrom Twitter", 2014, International Journal of Innovative Research in Computerand Communication Engineering, Vol. 2, Issue 12
19. Hassan Saif, YulanHe, Miriam Fernandez, and Harith Alani," Semantic Patterns for Sentiment Analysis of Twitter", 2014, ISWC Part II, LNCS 8797, pp. 324–340

20. SanthiChinthala, Ramesh Mande, SuneethaManne, and SindhuraVemuri," Sentiment Analysis on Twitter Streaming Data", 2015, Emerging ICT for Bridging the Future – Volume 1, pp- 470-481

21. Hassan Saif, Yulan He, and Harith Alani," Semantic Sentiment Analysis of Twitter", 2012, ISWC Part I, LNCS 7649, pp. 508–524

22. Syed Akib Anwar Hridoy, M. TahmidEkram, Mohammad Samiul Islam, Faysal Ahmed and Rashedur M. Rahman," Localized twitter opinion mining usingsentiment analysis", 2015, Anwar Hridoy et al. Decis. Analysis, vol. 4, issue 65 pp- 015-016

23. Xing Fang and Justin Zhan," Sentiment analysis using product review data", 2015, Springer, volume 5 issue 7, pp- 015-020

24. Khaled Ahmed, Neamat El Tazi, Ahmad HanyHossny," Sentiment Analysis Over Social Networks: AnOverview", 2015, IEEE, vol. 9, iss. 8, pp- 97-110

25. Aldo Hernández, Victor Sanchez, Gabriel Sánchez, Héctor Pérez,Jesús Olivares, Karina Toscano, Mariko Nakano and Victor Martinez," Security Attack Prediction Based onUser Sentiment Analysis of Twitter Data", 2016, IEEE, vol. 56, pp.45

26. Dan Cao, LiutongXu. Analysis of Complex Network Methods for Extractive Automatic Text Summarization.2016 2nd IEEE International Conference on Computer and Communications, vol. 9, iss. 8, pp- 97-110, 2016.

27. RasimAlguliyev, RamizAliguliyev, NijatIsazade. A Sentence Selection Model and HLO Algorithm for Extractive Text Summarization, IEEE, vol. 9, iss. 8, pp- 97-110, 2016.

28. NarendraAndhale, L.A. Bewoor. An Overview of Text Summarization Techniques. IEEE, vol. 9, iss. 8, pp- 97-110, 2016.

## AUTHORS PROFILE

**Prabhjot Kaur** is persuing Master of Engineering in Computer Science and Applications in Big Data specialization from Chandigarh University, Gharuan, India. She has done Bachelors of Engineering from Chitkara University, Baddi, Himachal Pradesh, India (2016). Her research interest includes machine learning and data analysis.

**Rajdavinder Singh Boparai** is Assistant Professor in Apex Institute of Technology of Chandigarh University, India. He has published more than twenty-five research papers in various journals and conferences. His primary research domain is data science and big data engineering and secondary research areas are Digital image processing and parallel computing.

**Dilbag Singh** received his PhD degree from Department of Computer Science and Engineering, Thapar Institute of Engineering and Technology, Patiala. He has done his Master in Technology (Computer Science and Engineering) from Guru Nanak Dev University, Amritsar, Punjab, India (2012). Currently, he is working as an assistant professor at Chandigarh University, Gharuan, Mohali, Punjab, India. He has published more than 27 research papers in well-known reputed SCI indexed journals and international conferences. His research interest includes Wireless sensor networks, Digital image processing and Meta-heuristic techniques.