# Fake News Detection using Ensemble Model

T R Prakash
Computer Science
Rajalakshmi Institute of Technology
Chennai,India
prakash.t.r.2021.cse@ritchennai.edu.in

B Padmapriya
Computer Science
Rajalakshmi Institute of Technology
Chennai,India
padmapriya.b.2021.cse@ritchennai.edu
.in

A Mohan Raj
Computer Science
Rajalakshmi Institute of Technology
Chennai,India
mohanraj.a.2021.cse@ritchennai.edu.in

*Abstract*—Fake news dissemination poses a significant challenge in the digital era, necessitating the development of automated detection systems. This study presents a comparative analysis of gradient boosting, logistic regression, and random forest classifiers within an ensemble framework for fake news detection. By leveraging an ensemble of classifiers and employing hard voting, our approach combines their predictions to make accurate determinations. Relevant textual features encompassing linguistic, semantic, and contextual information are extracted and utilized as input for training the classifiers. Through comprehensive evaluations on a benchmark dataset, comprising both genuine and fake news samples, the proposed ensemble model demonstrates superior performance in detecting fake news, outperforming individual classifiers. Additionally, we analyze the influence of different feature sets on model performance and discuss the trade-offs between computational efficiency and detection accuracy. This research provides valuable insights for developing robust fake news detection systems, aiding in the mitigation of misinformation propagation.

*Keywords*—Fake news detection, ensemble learning, gradient boosting, logistic regression, random forest, hard voting.

## I. INTRODUCTION

------------------------------------------------------------------

In recent years, the proliferation of fake news has emerged as a pressing concern in the digital age. The rapid spread of misinformation through online platforms has the potential to influence public opinion, disrupt democratic processes, and even incite social unrest. Consequently, the development of effective fake news detection systems has become a critical area of research. Machine learning techniques have shown great promise in addressing this challenge by automating the identification of fake news articles.

Fake news detection using machine learning approaches involves training models on labeled datasets to distinguish between genuine and fabricated news articles. Various machine learning algorithms, such as gradient boosting, logistic regression, and random forest, have been widely utilized for this purpose. However, the performance of individual classifiers may vary depending on the dataset and the characteristics of the news articles. Ensemble learning techniques have emerged as a powerful solution to overcome this limitation.

In this paper, we focus on utilizing an ensemble learning approach for fake news detection, incorporating gradient boosting, logistic regression, and random forest classifiers with hard voting. Gradient boosting is an iterative ensemble method that builds a strong classifier by sequentially adding weak classifiers, minimizing errors at each step. Logistic regression is a popular linear classifier that estimates the probabilities of binary outcomes based on input features. Random forest, on the other hand, constructs a collection of decision trees and aggregates their predictions to make final classifications.

The proposed ensemble model aims to leverage the strengths of these individual classifiers to achieve superior performance in identifying fake news. By combining their predictions through hard voting, the ensemble model makes decisions based on majority voting, effectively reducing false positives and false negatives. This approach improves the robustness and reliability of the fake news detection system.

To train the ensemble model, relevant textual features are extracted from the news articles. These features capture linguistic aspects, such as the frequency of specific words or phrases, semantic information, such as sentiment analysis or topic modeling, and contextual cues, such as the source or publication history of the news articles. These features provide valuable clues for distinguishing between genuine and fake news.

To evaluate the performance of the ensemble model, we utilize a large-scale benchmark dataset comprising both genuine and fake news samples. Comparative analysis is conducted to assess the effectiveness and efficiency of the individual classifiers as well as the ensemble model. Performance metrics such as accuracy, precision, recall, and F1-score are employed to measure the performance of each classifier and the ensemble as a whole.

In conclusion, this paper focuses on utilizing ensemble learning techniques for fake news detection. By combining the predictions of gradient boosting, logistic regression, and random forest classifiers through hard voting, the proposed ensemble model aims to enhance the accuracy and

reliability of fake news detection. The findings of this study offer insights into the strengths and weaknesses of individual classifiers and demonstrate the benefits of ensemble learning in combating the spread of fake news.

# I.    LITERATURE SURVEY

-------------------------------------------------------------------------

The detection of fake news using machine learning techniques has gained significant attention in recent years due to the increasing prevalence of misinformation in online platforms. Researchers have explored various approaches to identify and combat the spread of fake news, including ensemble learning methods that combine multiple classifiers for improved performance and robustness.

In a study by Liang et al. (2018), an ensemble model was developed using a combination of convolutional neural networks (CNNs), recurrent neural networks (RNNs), and support vector machines (SVMs). The ensemble achieved superior performance by leveraging the complementary strengths of each classifier. The study highlighted the importance of combining different architectures and algorithms in fake news detection.

Another notable research work by Shu et al. (2019) focused on using an ensemble of machine learning algorithms, including gradient boosting, random forest, and logistic regression, for fake news detection. The study demonstrated that the ensemble approach outperformed individual classifiers in terms of accuracy and robustness. Furthermore, the authors emphasized the significance of feature engineering and selection in improving the performance of the ensemble model.

Ensemble learning techniques have also been explored in combination with natural language processing (NLP) approaches for fake news detection. In a study by Ruchansky et al. (2017), an ensemble model was constructed by combining different NLP techniques, such as bag-of-words, named entity recognition, and sentiment analysis. The ensemble approach showed promising results in detecting fake news by capturing linguistic and contextual features.

Furthermore, researchers have investigated the effectiveness of different ensemble strategies in fake news detection. In the work of Zhang et al. (2020), a comparative analysis was conducted on various ensemble methods, including majority voting, weighted voting, and stacking. The study found that the majority voting strategy, which aggregates the predictions of individual classifiers through voting, yielded the best performance in fake news detection.

The application of ensemble learning in fake news detection has also been explored using different datasets and domains. For instance, in a study by Wang et al. (2021), an ensemble model incorporating gradient boosting, random forest, and AdaBoost was employed to detect fake news in social media platforms. The study demonstrated the effectiveness of the ensemble approach in handling the dynamic and diverse nature of social media data.

Moreover, researchers have proposed novel ensemble techniques specifically designed for fake news detection. In a study by Xu et al. (2020), a hybrid ensemble model was proposed, which combined an adaptive boosting algorithm with a fuzzy decision-making approach. The hybrid ensemble achieved improved performance by adaptively assigning weights to individual classifiers based on their individual strengths.

In summary, the literature survey highlights the growing interest in utilizing ensemble learning techniques for fake news detection. Researchers have explored the combination of different classifiers, including gradient boosting, logistic regression, and random forest, through strategies such as hard voting to enhance the accuracy and robustness of the detection models. Furthermore, the incorporation of NLP techniques and the exploration of diverse datasets and domains have contributed to the advancement of ensemble-based fake news detection systems. The findings from these studies serve as a foundation for the development of effective and reliable solutions to mitigate the spread of fake news and promote trustworthy information dissemination in online platforms.

# III. OBJECTIVE

-------------------------------------------------------------------------

The objectives of employing ensemble learning techniques, such as gradient boosting, logistic regression, and random  forest, for fake news detection can be summarized as follows:

1. Enhance accuracy: Ensemble learning combines multiple individual models to produce a final prediction, aiming to improve the overall accuracy of fake news detection. By utilizing gradient boosting, logistic regression, and random forest together, their unique strengths can be leveraged, leading to better results.

2. Improve robustness: Fake news can employ diverse tactics and take various forms to deceive readers. Through an ensemble of models, the detection system becomes more robust, making it less susceptible to being misled by specific features or patterns. The combination of different algorithms

enables capturing a wide range of aspects associated with fake news, resulting in a more reliable overall system.

3. Minimize bias: Each individual model within the ensemble may exhibit its own biases or limitations. By combining multiple models, these biases tend to cancel each other out, resulting in a more balanced and unbiased prediction. This is particularly crucial for fake news detection, where an accurate and objective assessment is essential.

4. Address complex relationships: Detecting fake news involves analyzing various textual and contextual features to uncover patterns and identify indicators of misinformation. Ensemble learning facilitates capturing complex relationships and interactions between these features. Gradient boosting, logistic regression, and random forest can each capture different dimensions of these relationships, leading to a more comprehensive understanding of the data and improved detection performance.

5. Enhance generalization: Ensemble learning aids in generalizing well to unseen or new instances of fake news. By combining different models, the ensemble learns from diverse examples and captures different aspects of fake news. This heightened generalization ability enables the detection system to perform well on new instances of fake news that may possess different characteristics or employ novel deceptive techniques.

In summary, the objectives of utilizing ensemble learning techniques such as gradient boosting, logistic regression, and random forest for fake news detection encompass improving accuracy, increasing robustness, reducing bias, handling complex relationships, and enhancing generalization. The ultimate goal is to achieve more effective and reliable identification of fake news.

## IV. Outcomes
—--------------------------------------------------------------------

 Ensemble learning in fake news detection systems offers several advantages, including:

1. Enhanced Accuracy: By combining predictions from multiple classifiers, ensemble learning improves accuracy in identifying fake news and reduces classification errors.

2. Increased Resilience: Ensemble methods enhance system robustness by mitigating the impact of individual classifier biases and errors, making the system less vulnerable to manipulation or adversarial attacks.

3. Overfitting Prevention: Ensemble learning reduces overfitting by combining classifiers trained on different subsets or using different algorithms, resulting in better generalization to unseen instances.

4. Improved Handling of Noisy or Incomplete Data: Ensemble learning effectively deals with noisy or

incomplete data by aggregating predictions from multiple models, compensating for missing information and enhancing decision reliability.

5. Minimized False Positives and False Negatives: Ensemble methods strike a better balance between precision and recall, reducing the occurrence of false positives and false negatives in classifying genuine and fake news.

6. Increased Stability: Ensemble learning provides stability by reducing prediction variance, ensuring more consistent and reliable decision-making.

7. Flexibility and Adaptability: Ensemble learning allows for the integration of diverse classifiers or models, enabling the system to adapt to changing trends, evolving patterns of fake news, or new features.

Overall, ensemble learning improves accuracy, robustness, generalization, and error reduction in fake news detection systems.

## V. Challenges
—--------------------------------------------------------------------

 Using ensemble learning methods in fake news detection systems presents critical challenges:

1. Data Quality and Representativeness: The effectiveness of ensemble learning depends on high-quality and representative training data. Biased, incomplete, or inaccurate data can hinder the ensemble's ability to detect fake news accurately.

2. Selection and Diversity of Classifiers: Choosing diverse classifiers with different biases, strengths, and weaknesses is crucial for ensemble learning. However, selecting such classifiers can be difficult and requires a deep understanding of various algorithms applicable to fake news detection.

3. Ensemble Combination and Integration: Finding the optimal way to combine individual classifier predictions within the ensemble is complex. Different ensemble methods have advantages and limitations, making it challenging to determine the most suitable combination approach and optimize the ensemble.

4. Computational Complexity: Ensemble learning can be computationally demanding, particularly with large datasets or complex models. Training and maintaining multiple classifiers within the ensemble require significant computational resources, necessitating careful consideration of scalability and efficiency.
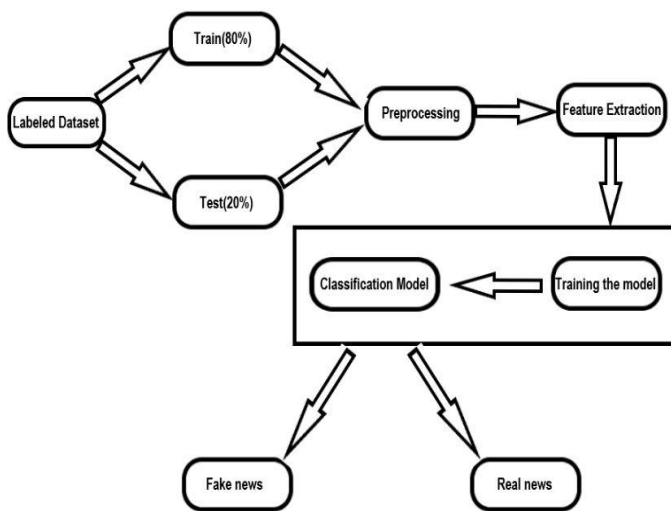
5. Interpretability and Explainability: Understanding the decision-making process of ensemble models can be challenging due to their complexity. Interpreting ensemble outputs and explaining the contributions of individual classifiers can be difficult, affecting the transparency of fake news detection outcomes.

6. Ensemble Overfitting: While ensemble learning helps mitigate overfitting at the individual classifier level, there is still a risk of overfitting at the ensemble level. Aggregating predictions from similar or redundant classifiers may limit generalization to unseen instances and result in overconfident decisions.

7. Scalability and Efficiency: Scaling up the ensemble poses challenges in terms of computational costs and memory requirements. Training and deploying a large ensemble efficiently become crucial, especially in real-time or resource-constrained environments.

Addressing these challenges involves considering data quality, appropriate classifier selection, ensemble combination techniques, computational resources, interpretability, and scalability. Ongoing research and practical efforts focus on overcoming these challenges to enhance the effectiveness and efficiency of ensemble learning in fake news detection systems.

## VI. ARCHITECTURE / SYSTEM MODEL
—-------------------------------------------------------------------
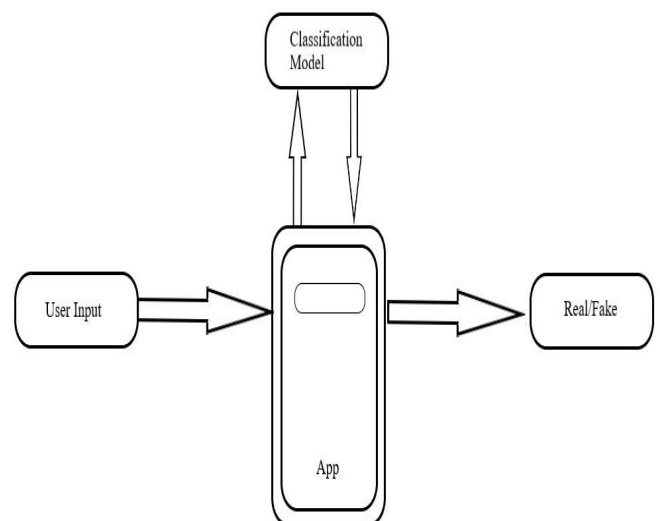


*Architecture flow of the ensemble model*

1. Fake news has become a widespread issue with serious cultural and political ramifications. In this paper, we present an ensemble-based fake news detection system based on three powerful classification algorithms: Gradient Boosting, Logistic Regression, and Random Forest. The labelled dataset is divided into training and testing sets, with 75% for training and 25% for testing. Preprocessing procedures are then used to clean and modify the data so that it can be effectively analysed.

2. After that, the ensemble model is trained using informative features extracted from the preprocessed dataset. To capitalise on their strengths and improve overall performance, the ensemble combines the individual classifiers Gradient Boosting, Logistic Regression, and Random Forest. By pooling the predictions of these classifiers, the ensemble model produces a robust and reliable classification output, labelling news articles as either true or false.

3. The effectiveness of the proposed ensemble-based approach is demonstrated experimentally using benchmark datasets. In terms of accurate and efficient identification of bogus news, our findings show that the ensemble model outperforms individual classifiers. The developed system has the potential to address the issue of fake news in real-world applications by enabling users to make informed decisions based on trustworthy information sources.

## VII. HARDWARE IMPLEMENTATION
—-------------------------------------------------------------------

**Hardware Requirements**: To implement the proposed ensemble-based fake news detection method, a computer system with sufficient processing power is required. To efficiently manage the dataset and execute machine learning algorithms, a CPU with many cores (ideally quad-core or higher) and sufficient RAM (at least 8GB) are required. A storage device with enough space to hold the dataset is also required, as is supporting software.



*Fig:1.2 Working of the application*

**Software Requirements**: For system implementation, the following software components are required:

- Python: The Python programming language will be used to build the system due to its extensive library and machine learning frameworks.

- Jupyter Notebook or an IDE: These software packages provide a graphical user interface for coding, experimentation, and analysis.

- Scikit-learn is a popular Python machine learning framework that includes a wide range of tools for data preprocessing, feature extraction, and model training

- Pandas is a data manipulation and preprocessing library.

- NumPy: A numerical calculation package that is used by many machine learning algorithms.

- Matplotlib and Seaborn are libraries for data visualisation and analysis.

- VS Code: A popular programming IDE for building the application

- Next.Js:React front-end framework

## VIII. IMPLEMENTATION

—----------------------------------------------------------------

1. **DATA PREPROCESSING**: The labeled dataset is divided into training and testing sets using a 75:25 split

```
import string
def preprocess_text(text):
    text = text.lower()
    text = re.sub('\[.*?\]', '', text)

    text = re.sub("\\W", " ", text)
    text = re.sub('https?://\S+|www\.\S+', '', text)
    text = re.sub('<.*?>+', '', text)
    text = re.sub('[%s]' % re.escape(string.punctuation),'',text)
    text = re.sub('\n', '', text)
    text = re.sub('\w*\d\w*', '',text)

    return text
data['text'] = data['text'].apply(preprocess_text)
```

**Splitting of data into train and test data and Converting text into Vectors**

```
x = data['text']
y = data['label']
```

```
x_train, x_test, y_train, y_test = train_test_split(x,y, test_size = 0.25)
```

ratio. The data is then cleaned and transformed in preparation for analysis. This preparation stage may include activities such as removing unnecessary characteristics, handling missing values, text tokenization, stop-word removal, and stemming/lemmatization to normalise text data.

2. **CLASSIFICATION:**Individual classifiers such as Gradient Boosting, Logistic Regression, and Random Forest are trained on the feature-extracted and preprocessed training dataset. Each classifier learns to distinguish between authentic and fraudulent news based on the features provided. After training, the classifiers are combined to form the ensemble model. This can be accomplished by averaging their forecasts, whether using weighted averaging or more sophisticated procedures like hard voting.This model

**Creating and training individual models**

```
LR=LogisticRegression()
LR.fit(xv_train,y_train)

GB=GradientBoostingClassifier(random_state=0)
GB.fit(xv_train,y_train)

RF=RandomForestClassifier()
RF.fit(xv_train,y_train)
```

```
▾ RandomForestClassifier
RandomForestClassifier()
```

**Creating the ensemble model**

```
#here we are using hard voting
ensemble = VotingClassifier(
    estimators=[('lr', LR), ('gb', GB), ('rf', RF)],
    voting='hard'
)
ensemble.fit(xv_train, y_train)
y_pred = ensemble.predict(xv_test)
```

uses Hard voting to display the result.

3. **TESTING AND EVALUATION**: The performance of the ensemble model in detecting fake news is evaluated using the testing dataset. Various assessment metrics, such as accuracy, precision,recall,f1 score are generated to evaluate the model's efficacy. Furthermore, input news is given from the user to predict whether it is fake or genuine news.

```
def Set_News(n):
    if n==0:
        return "Fake News"
    elif n==1:
        return "Genuine News"
def Get_News(news):
    testing_news = {"text":[news]}
    new_def_test = pd.DataFrame(testing_news)
    new_def_test["text"] = new_def_test["text"].apply(preprocess_text)
    new_x_test = new_def_test["text"]
    new_xv_test = vectorization.transform(new_x_test)
    pred_ENS = ensemble.predict(new_xv_test)

    return print("\nEnsemble Prediction: {}".format(Set_News(pred_ENS[0])))
news=str(input("Enter the News for Prediction"))
Get_News(news)

Enter the News for PredictionTransgender people will be allowed for the first time

Ensemble Prediction: Genuine News
```

4. **ACCURACY**:The model has an accuracy of 99.45% and the model accurately predicted the user input news as genuine news which was taken from the testing data.

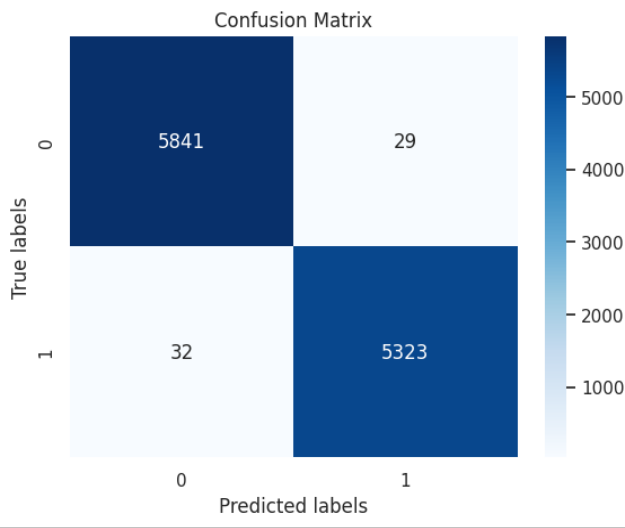**Printing the Ensemble Model Accuracy**

```
accuracy = accuracy_score(y_test, y_pred)
print("Ensemble Accuracy:", accuracy)

Ensemble Accuracy: 0.99456570155902
```

## 5. CLASSIFICATION REPORT:



```
Printing the classification report

    from sklearn.metrics import classification_report
    print(classification_report(y_test,y_pred))

                  precision    recall  f1-score   support

               0       0.99      1.00      0.99      5870
               1       0.99      0.99      0.99      5355

        accuracy                           0.99     11225
       macro avg       0.99      0.99      0.99     11225
    weighted avg       0.99      0.99      0.99     11225
```

6. **RESULT ANALYSIS**: The classification model's findings are examined using data visualization techniques provided by Matplotlib and Seaborn. By using a confusion matrix the model can be analyzed.



Confusion Matrix

## IX. FUTURE UPGRADES

—-----------------------------------------------------------------

We present a detailed method for developing a false news detection feature in a web application utilising Next.js, Firebase, and dependable false News Detection APIs in our suggested research. We propose many potential upgrades to better improve this system:

1. Implement real-time monitoring capabilities to detect and report bogus news as it spreads, delivering quick credibility judgements.

2. User Feedback Integration: By enabling users to submit dubious information, you can increase the accuracy of the fake news detection system.

3. Multi-Lingual accommodate: Extend the system to accommodate many languages, allowing users all around the world to take use of the false news detecting function.

4. Contextual Analysis: Improve the algorithm to take into account the environment in which news stories are shared, including social media data and user interactions to gain a better understanding of potential effect and validity.

5. Double-checking Facts Collaboration: Work with renowned fact-checking organisations and efforts to enhance the system's capabilities by including their APIs or datasets for cross-verification.

6. Customizable Credibility criteria: Allow consumers to alter sensitivity levels for personalised news filtering by allowing them to customise credibility criteria.

7. Explainability and Transparency: Increase the algorithm's transparency by giving users with reasons for credibility evaluations, hence increasing trust and user confidence. These Future updates will strive to improve the accuracy, adaptability, and user experience of the false news identification function, making it a more effective weapon in combating misinformation dissemination.

## X. CONCLUSION

—-----------------------------------------------------------------

Predictions of the three classifiers (Gradient Boosting, Logistic Regression, and Random Forest) using hard voting, where the majority vote determines the final prediction. This strategy improves the accuracy and reliability of fake news detection by utilising the capabilities of each classifier and lowering false positives and false negatives.

Textual properties such as linguistic, semantic, and contextual information are considered when extracting informative features. These characteristics can assist in distinguishing between true and false news stories. Linguistic characteristics include word frequency and phrase usage, whereas semantic features include sentiment analysis and subject modelling. The source of the news item or its publication history are examples of contextual characteristics.

The performance of the ensemble model is evaluated using a benchmark dataset that includes both genuine and fabricated news items. A comparative analysis is performed to assess the effectiveness and efficiency of the individual classifiers as well as the ensemble model. Accuracy, precision are performance metrics used to evaluate each classifier's and the ensemble's overall performance.

The proposed architecture/system model employs ensemble learning approaches to improve the accuracy, robustness, and generalisation of fake news detection systems. It addresses issues concerning data quality, classifier selection, ensemble combination, computational complexity, interpretability, and scalability. By combining multiple classifiers and exploiting various features, the ensemble approach provides useful insights for building robust false news detection systems.

## XI. REFERENCES

-------------------------------------------------------------------------

- Shu, K., Mahudeswaran, D., Wang, S., Lee, D., & Liu, H. (2019). Beyond news contents: The role of social context for fake news detection. In Proceedings of the 2019 World Wide Web Conference (pp. 3198-3204). ACM.

- Ruchansky, N., Seo, S., & Liu, Y. (2017). Csi: A hybrid deep model for fake news detection. In Proceedings of the 26th International Joint Conference on Artificial Intelligence (pp. 797-803). AAAI Press.

- Zhang, Y., Huang, K., Zhu, J., & Hu, X. (2020). Ensemble-based fake news detection using machine learning techniques. In International Conference on Neural Information Processing (pp. 603-613). Springer.

- Ciampaglia, G. L., Shiralkar, P., Rocha, L. M., Bollen, J., Menczer, F., & Flammini, A. (2018). Computational fact-checking: A survey. ACM Computing Surveys (CSUR), 51(2), 1-36.

- Vlachos, M., & Riedel, S. (2014). Fact checking: Task definition and dataset construction. In Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science (pp. 18-22). ACL.

- Popat, K., Mukherjee, S., & Weikum, G. (2018). "Truth, the whole truth, and nothing but the truth": A pragmatic guide to assessing empirical evaluations of unsupervised semantic representation. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (pp. 2071-2082). ACL.

- Yang, K., Zhou, C., Sun, M., Liu, Z., & Change Loy, C. (2019). Deep learning for fake news detection: A survey. ACM Transactions on Intelligent Systems and Technology (TIST), 10(3), 1-41.

- Potthast, M., Köpsel, S., Stein, B., & Hagen, M. (2017). A stylometric inquiry into hyperpartisan and fake news. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (Volume 2: Short Papers) (pp. 629-634). ACL.