

Detecting Fake News Over Job Posts via Bi-Directional Long Short-Term Memory (BIDLSTM)

T. V. Divya, Koneru Lakshmaiah Education Foundation (Deemed), Hyderabad, India

 <https://orcid.org/0000-0003-0167-9599>

Barnali Gupta Banik, Koneru Lakshmaiah Education Foundation (Deemed), Hyderabad, India

 <https://orcid.org/0000-0002-8107-9501>

ABSTRACT

Detection of fake news on job advertisements has grabbed the attention of many researchers over the past decade. Various forms of classifiers such as support vector machine (SVM), XGBoost classifier, and random forest classifier (RF) methods are greatly utilized for fake and real news detection about job advertisement posts in social media. There exhibits the slight or elusive variance among fake and the real news, which are obtained through topics and word embeddings, that affect system accuracy. Initially, pre-processing steps for job post data like stop word removal, tokenization, and lemmatizing words are done utilizing wordnet. The oversampling procedure is accomplished for data balancing. Subsequently, the new columns are generated representing each possible attribute value existence from original data by suggesting one-hot encoding. The dataset insignificant features removal is accomplished, which is exploited for fake news detection. As a final point, bi-directional long short-term memory classifier (Bi-LSTM) is greatly utilized for learning word representations in lower-dimensional vector space and learning significant words word embedding or terms revealed through the word embedding algorithm. The fake news detection is greatly achieved along with real news on job posts from online social media which is achieved by Bi-LSTM classifier, thereby evaluating corresponding performance. The performance level metrics such as precision, recall, F1-score, and accuracy are assessed for effectiveness by fraudulency based on job posts. The outcome infers the effectiveness and prominence of features for detecting false news.

KEYWORDS

Bi-Directional Long Short-Term Memory (Bi-LSTM), Fake News Detection, Job Advertisement Posts, Machine Learning, Social Media, Text Processing

INTRODUCTION

As stated by the United States (US) Department of Labor, the rate of unemployment is 11.1% in the Bureau of Labor Statistics US Department of Labor, Employment Situation of US as of June 2020., Even though a lot of factors exist behind the present unemployment rate, several people are there in the US as well as in other parts of the world, those who look forward to getting new jobs because of the job loss and some other financial crisis. In recent days, almost all companies have enabled posting

DOI: 10.4018/IJWLTT.287096

This article published as an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

in job boards directly or pulling job data from the job aggregators, since most of the job postings have been performed online. Nevertheless, it is unsure that every job postings are real, since some of them are fraudulent that may be intended to obtain data or other confidential details from desperate job seekers . It has been posted in USC Career Center. Avoid Fraudulent Job Postings. Accessed July 2020.

According to Alghamdi, B., & Alharby, F. (2019), recently in the area of Online Recruitment Frauds (ORF), an employment scam becomes one of the crucial problems. Even though it possesses the benefits (i.e. quick and easy access for the job seekers) and preferred by most organizations, yet performing a job posting on social media has considered being a double-edged sword. The reason behind all of these scams created by the fraudsters might be the cheating of money from the job-seekers by offering employment to them. Hence, it is mandatory to make sure about the job advertisement whether it has been posted on the real site or the fake one.

According to Ibrishimova, M. D., & Li, K. F. (2019) and Zhou, X., Zafarani, R., Shu, K., & Liu, H. It is the biggest challenge and highly impossible to identify these kinds of frauds. For conquering these challenges, the frameworks belong to Artificial Intelligence (AI) and Natural Language Processing (NLP) can predominantly get utilized as they facilitate the researchers to develop the classifier that makes automatic detection of fake news possible . It has considered being a crucial process at this time since numerous job-seekers apply online for jobs because of the present unemployment scenario. Besides, it is mandatory to escaping the scam victims from fake job postings to recover the economy. As described by Khan, J.Y., Khondaker, M., Islam, T., Iqbal, A. and Afroz, S., (2019). Next to NLP, The Machine Learning (ML) method has vitally utilized, since it applies various classification algorithms to identify the fake posts. During that, the fake job posts have segregated from a huge set of job ads with the help of the classification system and warn the user.

In previous studies, According to Ahmed, H., Traore, I., & Saad, S. (2017), Kaliyar, R. K. (2018) various conventional ML and Neural Networks (NNs) approaches have been exploited for detecting fake news . Nevertheless, As described by Singhanian, S., Fernandez, N., & Rao, S. (2017) those methods have solely concentrated on social and political news . Therefore, the frameworks were also developed according to their area of interest, since the features of the models have been designed for specific datasets. Consequently, they get trapped into dataset bias and deliver a poor performance over the news of some other topic. Nonetheless, it does not apply many of the enhanced machine learning approaches, such as neural network-based approaches, those which have ascertained to be the optimal solution in several text classification issues. Moreover, the previous studies significantly include the flaw that they have implemented over a specific kind of dataset. Due to this reason, the performance evaluation of several models gets intricate.

Moreover, in these works, the difference between the topics and word embeddings displays slight or refined modification between the fake and real news. It limits the prototype's capability to comprehend how far extend related or unrelated the reported news appears when compared to the original news. It reduces the accuracy of the system. In addition to depend completely on language, the method relies on remote n-grams, often removed from the suitable context info. Word embedding systems are mostly providing an useful way to represent the meaning of the word. At the same time existing works, a label encoding function is used for converting the text data into the numerical format. New encoding methods are required for labeling text data numerically in an understandable manner. In some kind of circumstances, sentences of diverse lengths could be signified as a tensor with altered dimensions. Traditional models cannot handle the sparse and high order topographies quite well.

The Deep Learning (DL) models are most commonly used in both the academic community and industry. In the NLP area, DL models are deployed to train a model that could denote words as vectors. According to Chen, K., Wang, J., Chen, L. C., Gao, H., Xu, W., & Nevatia, R. (2015), several researchers initiated numerous deep learning models based on the word vectors for Question Answering (QA), etc. Convolutional Neural Networks (CNN) utilizes filters to internment the local structures of the image, which performs very well on computer vision tasks. Bi-directional Long Short Term Memory (Bi-LSTM) performed better than the classifiers based on CNN.

In this article, the Bi-LSTM framework is proposed which detects and classifies fake job postings. First, visualize the insights from the fake and real job posts, then apply some pre-processing steps for detecting the data as fake or real. Thirdly attributes that are not used for detection are removed from the text pre-processed dataset. Before that one hot encoding is introduced for the text encoding and Word embedding is introduced for mapping from discrete objects such as words to vectors. Embedded vector is given as input for the prediction of the real and fraudulent class labels for the job advertisements after successful training. It can be trained on the previous real and fake job and it can identify a fake job accurately. Finally, evaluate the performance of classifiers using several evaluation metrics. This classifier is trained on a benchmark dataset collected from Kaggle to identify the fake job posts.

LITERATURE REVIEW

Machine Learning (ML) models tend to possess several use cases. Text Analytics (TA), a type of Natural Language Processing (NLP) permits ML algorithms utilized for classification models textual data. The outline of various ML approaches utilized for fake and the real news detection is discussed in this section.

According to Agarwal, V., Sultana, H. P., Malhotra, S., & Sarkar, A. (2019) utilized NLP and ML approaches for fake and the real news detection. The count vectorizer, bag-of-words, n-grams are used along with Term Frequency-Inverse Document Frequency(TF-IDF) and training of data on five classifiers such as Logistic Regression (LR), Support Vector Machine (SVM), Naïve Bayes(NB), Stochastic Gradient (SG) and Random Forest(RF) for investigation of labeled news statements with the specific dataset. Also, precision, recall, and F1-scores are examined for model analysis.

Poddar, K., & Umadevi, K. S. (2019) exploited probabilistic computational models and geometric machine learning models for mitigating fake news. The proper vectorizer for fake news detection is obtained by comparing two different vectorizer scores namely count and TF-IDF. The score improvement is achieved by English stop words. The fake news prediction is done through different classifiers such as NB, SVM, LR, and Decision Tree (DT) classifier. SVM with TF-IDF outclasses best results which are validated by simulation.

According to Mahabub, A. (2020) developed an Ensemble Voting Classifier system for performing both real and fake tasks in news classification. Familiar machine-learning systems like NB, SVM, RF, Artificial Neural Network (ANN), LR, Ada Boosting, K Nearest Neighbor (KNN) are utilized in this research for detection. The cross-validation is done for obtaining the best three ML algorithms utilized in the ensemble voting Classifier. It is revealed by the experimental outcomes that 94.5% accuracy is attained by the proposed framework along with improved Receiver Operating Characteristic curve (ROC) score, precision, recall, and F1-score. The news most significant highlights are also obtained by this methodology. It may be further utilized for fake profiles, fake message detection additionally.

Thota, A., Tilak, P., Ahluwalia, S., & Lohia, N. (2018) suggested Deep Learning (DL) architectures for fake news detection. An instantaneous necessity for inevitably tagging besides perceiving such perverse news articles is presented by the exponential increase in production and inaccurate news distribution. It also necessitates a model for understanding nuances in natural language which adds complexity for fake news automated detection. It is always considered as a binary classification task since major existing fake news detection models are handled similarly. This also restricts the model's ability for understanding related or unrelated reported news is associated with real news. Neural network architecture is used for bridging the gap for accurate prediction of stance amid a specified pair of headline and article body.

Zhang, J., Dong, B., & Philip, S. Y. (2020) presented an gated graph neural network called FAKEDETECTOR for fake news prediction based on explicit and latent features set mined from the textual information. FAKEDETECTOR creates a deep diffusive network model for learning news articles, creators, and subject representations instantaneously. FAKEDETECTOR is compared with traditional methods for real-world fake news datasets using wide-ranging experiments.

Yang, Y., Zheng, L., Zhang, J., Cui, Q., Li, Z., & Yu, P. S. (2018) offered a review on “fake news detection”. Fake news detection in an automatic manner is considered to be a challenging task due to model-based fact-checking for news. The valuable clear topographies were recognized massively from equally text words as well as images exploited in fake news after fake news data exhaustive examination. There persist few hidden patterns in words as well as images used in fake news apart from explicit features, which might be taken with a set of latent structures obtained by multiple convolutional layers in the TI-CNN model. This typical model exploits explicit and latent structures into an unified feature space and is also accomplished with both text and image information concurrently. The effectiveness of this method outperforms the traditional method which is revealed through extensive experimentation.

Monti, F., Frasca, F., Eynard, D., Mannion, D., & Bronstein, M. M. (2019) utilized geometric deep learning approach for fake news detection in an automatic manner. The traditional CNN has generalized to graphs by the fundamental algorithms, which enables the fusion of heterogeneous data, namely content, user profile and activity, social graph and news propagation. Based on the news articles, testing and training have been performed for the CNN framework. Besides, it has been validated through professional fact-checking organizations by circulating it on Twitter. Consequently, it can be ensured that the initial detection of fake news, within a short span after propagation. Next, the aging of the CNN system has been tested overtraining and testing data separated in time. The outcomes ensure that the propagation-based methodologies are the optimal replacement approach for content-based techniques in terms of fake news detection. Besides, the experiments depict the effectiveness of the social network structure and propagation to enable extremely accurate fake news detection 92.7% ROC.

Zhang, J., Cui, L., Fu, Y., & Gouza, F. B. (2018) focused on the reduction of difficulties that occurred by the unfamiliar attributes of fake news as well as multiple correlations that existed in news articles, subjects, and creators. Following the comprehensive investigation, a new automatic fake news credibility inference framework has been presented, namely FakeDetector. A deep diffusive network framework has been built by FakeDetector built on a set of obvious and latent features derived from the textual information, through which the representations of news articles, subjects, and creators have concurrently learned. Based on a real-world fake news dataset, the proposed system has comprehensively experimented with. Subsequently, the results obtained by the FakeDetector and other existing methods have been compared, through which the proposed approach has proved its efficiency.

Liu, Y., & Wu, Y. F. (2018) tend to perceive fake news on social media in the initial stage by classifying the news propagation paths, for which they proposed a new framework. At first, each news article’s propagation path has taken as a multivariate time series, where all tuples are numerical vector that indicates the attributes of a user who disseminates the news. Subsequently, by integrating both convolutional and recurrent networks, a time series classifier has been constructed. In fake news detection, this model efficiently captures the global and local variations of user attributes together with the propagation path, correspondingly. About three real-world datasets, the empirical outcomes depict the efficiency of the proposed method to secure 85% accuracy on Twitter, and 92% accuracy on Sina Weibo within 5 minutes, once it begins spreading. This performance has proved to be higher than other prevailing baselines.

Roy, A., Basak, K., Ekbal, A., & Bhattacharyya, P. (2018) tend to classify fake news into the pre-defined fine-grained categories, for which they designed several deep learning frameworks. Accordingly, CNN and Bi-LSTM networks have proposed initially fake news detection. Subsequently, the derived outcomes have given into a Multi-Layer Perceptron (MLP) to implementing the final classification. Through the empirical findings on a benchmark dataset, the optimal results have been obtained with a 44.87% of accuracy rate, which is superior to other innovative models.

Dutta, S., & Bandyopadhyay, S. K. (2020) Machine Learning-based classification methods, like RF, DT, NB, KNN, MLP, Gradient Boost Classifiers, and AdaBoost Classifiers, through which they proposed an automated framework. Web-based fraudulent post has investigated through various classifiers, and the optimal employment scam identification framework has recognized by comparing

the obtained results of those classifiers. It has the potential to carry out the fake job post-detection from the huge quantity of posts. Accordingly, a single classifier and ensemble classifiers have been taken into consideration during the process. Compared to the single classifier, the higher classification efficiency of ensemble classifiers has been proved through the empirical outcomes, concerning scam detection.

PROPOSED METHODOLOGY

Social media makes possible the propagation of ‘fake news’ (poor quality news that purposefully includes false information). Throughout this work, for detecting the job posts-related fake news on the social media, Deep Learning algorithm has presented. Accordingly, from the benchmark site, the dataset of job posts has been gathered. Post-collection of the job posts dataset carried out from benchmark site, a few steps need to be performed to classify the job posts as fake or real. The steps include Text Preprocessing, Feature Extraction, Data Oversampling, Data Encoding, Word Embedding, and Fake News Detection. In-Text Pre-processing step, further sub-steps have involved, i.e. knowing missing values, visual representation of text data (Word Cloud), Fill null values, Cleaning Text Features by Stop Words Removal, Tokenization, and Lemmatizing using Wordnet. During the step of Feature Representation, the current features get split into four categories, namely texture, complex, binary and categorical. Subsequently, for resolving the data imbalance problem, the oversampling approach has been utilized. Thereby the balanced dataset has been obtained, from which the text encoding has been carried out through the one-hot encoding method to eliminate the irrelevant attributes from the dataset. Subsequently, Word embedding has been further employed to map from discrete objects such as words to vectors and real numbers. Ultimately, fake news detection has accomplished by employing the Bi-Directional Long Short Term Memory (Bi-LSTM). Figure 1 illustrates complete architecture of the framework.

1.1. Fake Job Postings Dataset

Implementation and testing of the proposed approach have been carried out through predicting-fraudulency-based-on-job-posts”, which has made up of the dataset of job descriptions and associated meta-information. Through the column “fraudulent”, a small proportion of aforementioned descriptions has identified as fake/scam.

1.2. Visualizing Missing Values

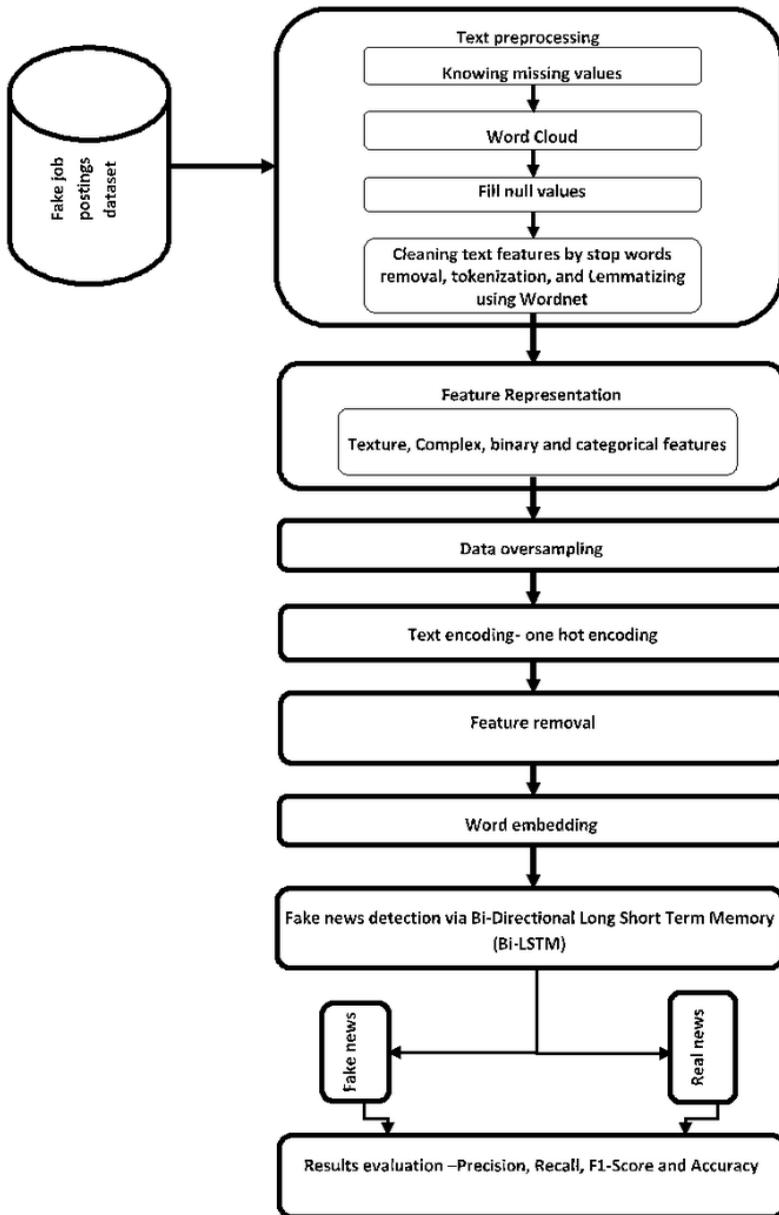
In the dataset, presented missing values have been visualized through the graphical framework, which has represented in Figure 2 and Figure 3 whereas, Figure 4 demonstrates the sample count, which is missed for individual features accompanied by their specific ranges.

Table 1 numerically represents the data with four features, i.e. job_id, telecommuting, has_company_logo, and has_questions. In the text classification, the data does not prove its significance, thus it has evaded from the dataset. Finally, the text data has visually represented by utilizing the remaining 14 features.

1.3. Wordcloud On Job Titles

Wordcloud is an approach that visually represents the text data, in which a list of words will be exhibited. Among that, the significance of each word will be highlighted through the color/size of the font. Through this format, the more frequent words can get perceived more quickly.

Figure 1. Proposed fake news detection system for job postings



FILL NULL VALUES

From Text Features, the Null Values has filled using the function called 'Unspecified', during which two features have predominantly utilized, namely salary range and the locations. Following that, the locations feature have further distributed into three kinds, i.e. Country, State, and City. Whereas, the salary feature has reformed into ranges within Min & Max by filling null values from salary_range with '0-0', during which the salary range has divided into two individual columns, i.e. min_salary and max_salary. Due to the different number of job posts in each country, the location feature has

Figure 2. Visualizing missing values

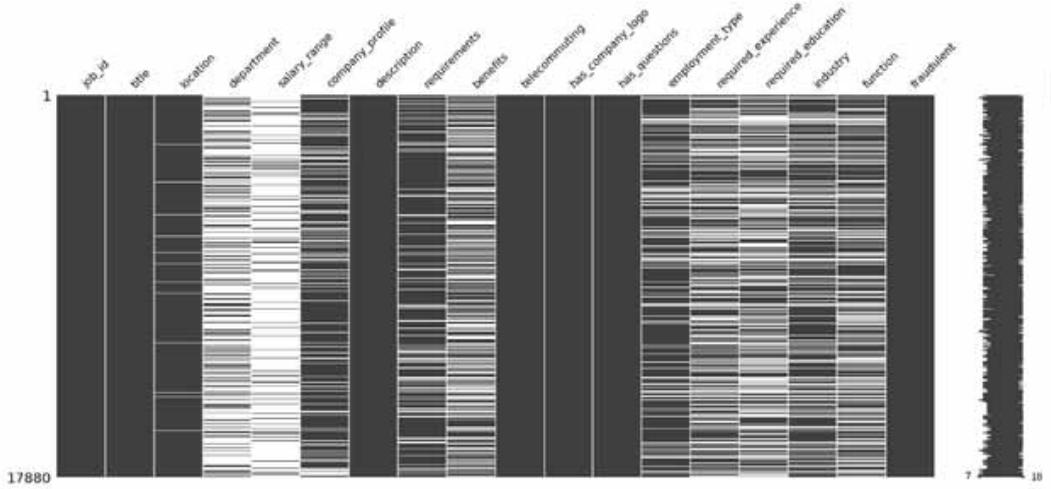


Figure 3. Visualizing missing values graph

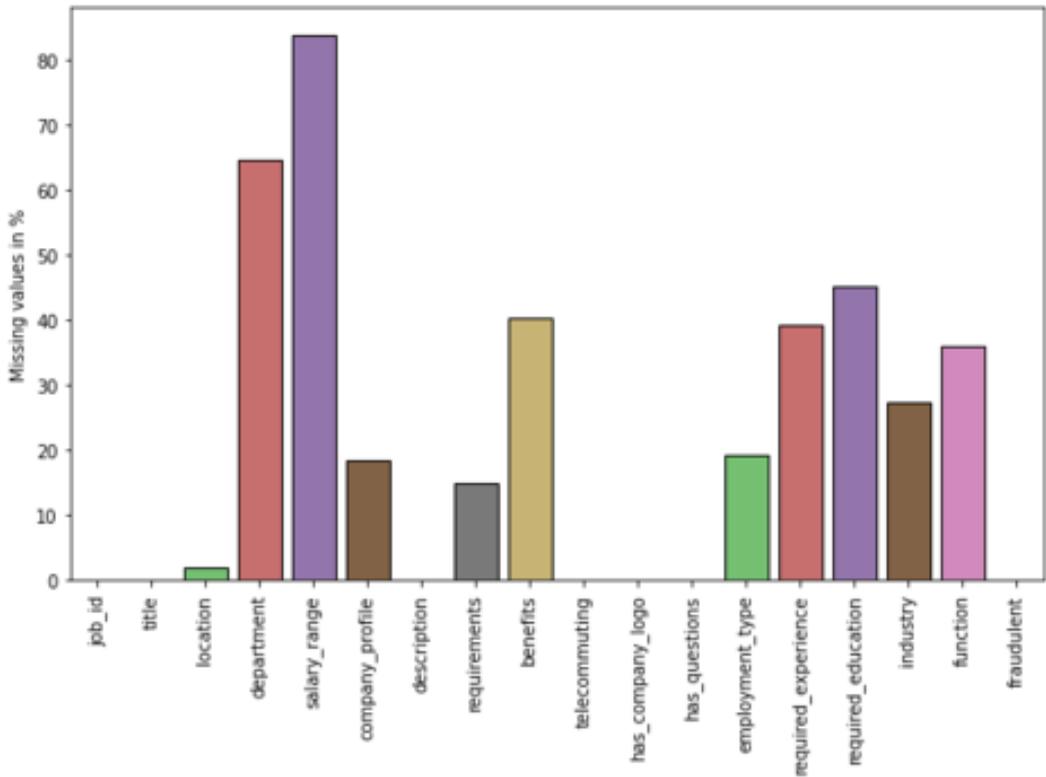


Figure 4. Missing values with their ranges

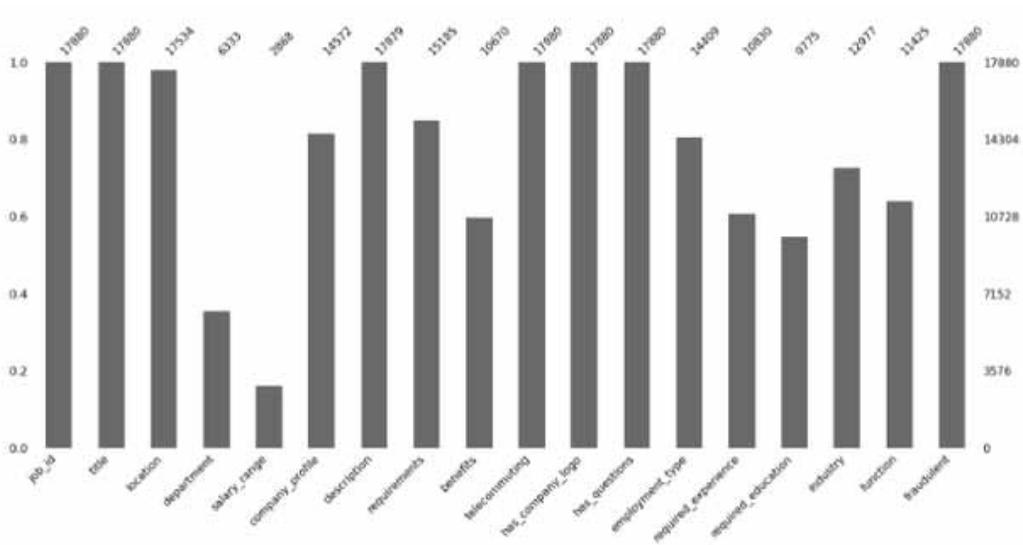


Table 1. Features removal with parameters

PARAMETER	job_id	telecommuting	has_company_logo	has_questions	Fraudulent
Count	17880.000000	17880.000000	17880.000000	17880.000000	17880.000000
Mean	8940.500000	0.042897	0.795302	0.491723	0.048434
Std	5161.655742	0.202631	0.403492	0.499945	0.214688
min	1.000000	0.000000	0.000000	0.000000	0.000000
25%	4470.750000	0.000000	1.000000	0.000000	0.000000
50%	8940.500000	0.000000	1.000000	0.000000	0.000000
75%	13410.250000	0.000000	1.000000	1.000000	0.000000
max	17880.000000	1.000000	1.000000	1.000000	1.000000

considered to be a significant metric for job postings. Figure 5 depicts the number of jobs posted country-wise.

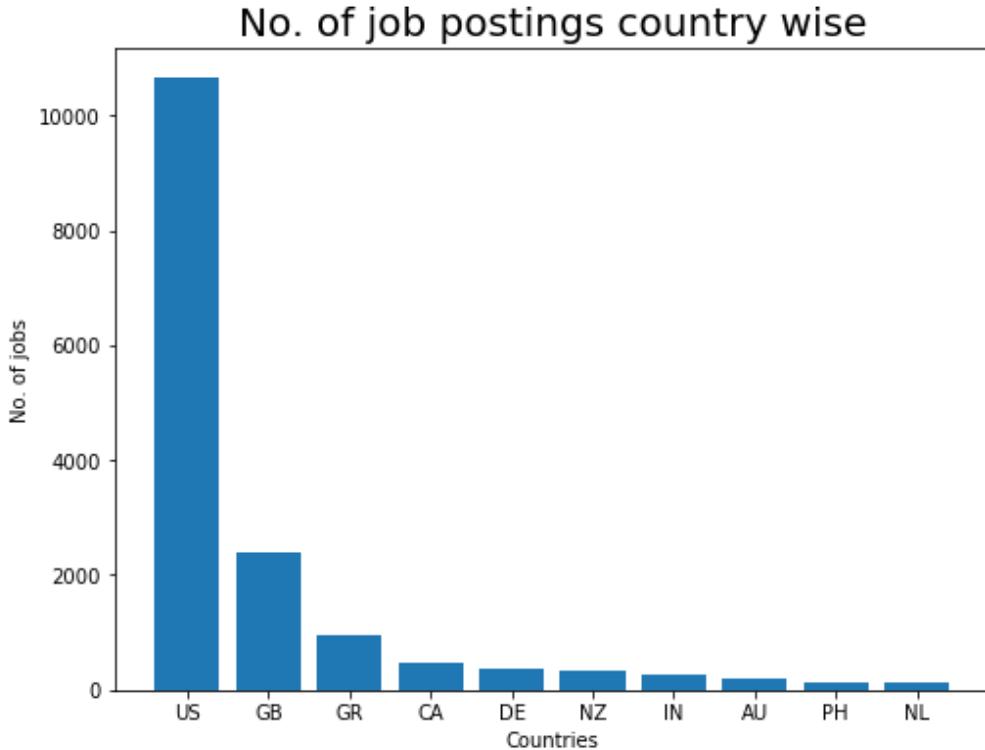
1.4. Cleaning Text Features

Cleaning of text features has primarily been carried out through stop words removal, tokenization, NLP helps in lemmatizing of words. For receiving column as an argument, applying regular expression functions, removing stop words, tokenizing, lemmatizing, and returning the modified data frame, the clean_text function has applied.

Argument: It has known to be the value. If an argument is called, it has sent to the function.
defclean_text(data)

Regular Expression: It has considered being a special sequence of characters. It aids in matching/ identifying other sets of strings or individual strings via specialized syntax held in a pattern,
 description = re.sub("[^a-zA-Z]", "", description)

Figure 5. Number of job postings in country wise



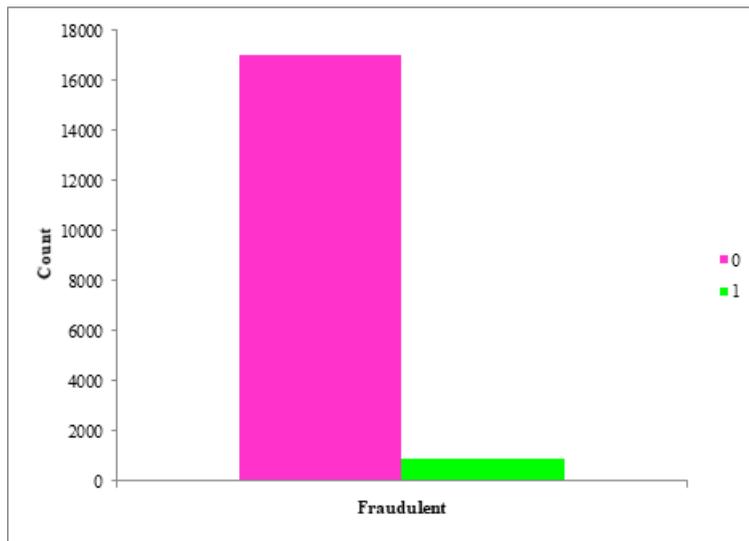
Stopwords Removal: Stopwords refer to the English word that does not make much sense in the sentence, and they could be omitted with caution, regardless of compromising the sense of the sentence (e.g. he, have, the, etc.). Such words are already captured in a corpus named corpus. Subsequently, generated a stop words list, and it has compared with the source document. Post-comparison, if any matching words are found in this list that will get eliminated from the source file. Besides, stop words convey no value in equivalent queries to the document. Hence, the removal of stop words not impacts the semantics of the document.

Millstein, F. (2020), Hardeniya, N., Perkins, J., Chopra, D., Joshi, N., & Mathur, I. (2016) has proposed the **Tokenization:** The task, where the text has segmented into words, clauses/sentences is called Tokenization. During this progress, bulk texts will get divided into tinier parts that have termed as tokens which help to identify such patterns and serves as a base step for lemmatization. Here, the approach, namely word_tokenize() has applied to divide a sentence into words. Following that, to optimally understand the text in the machine learning applications, the word tokenization output can get converted to a Data frame.

description=nlp.word_tokenize(description) is used for tokenization.

Bird, S., Klein, E., & Loper, E. (2009), Perkins, J. (2010) proposed **Lemmatization:** During this process, various inflected forms of the word have organized together that eases the analysis, as they have assessed as a single item. The Lemmatization process seems identical to stemming, but it sets the context to the words. Even though Lemmatization has a close correlation with stemming, yet stemming gets distinguished by operating on a single word deprived of the context knowledge, thereby it is inadequate to discriminate the words that have different meanings which are subject to part of speech. Wordnet has known to be an open-access and huge lexical database that tends to create structured semantic relationships within words, particularly in the English language. Besides,

Figure 6. The distribution of the target feature (fraudulent)



this database enables the potentials for performing the lemmatization. And, it has considered being one of the eldest and widely-recognized lemmatizers, for which Natural Language Toolkit (NLTK) facilitates an interface .

1.5. Oversampling Target Variable

Post-completion of the preprocessing task, the fake news detection can get initiated. As represented in Figure 6, among overall job posts, Target Variable-fraudulent from 866 jobs are fraud, whereas 17014 jobs are real. Consequently, the data-oversampling process has been carried out since the detection process turns out to be challenging to execute.

1.6. Text Encoding- One-Hot Encoding

In deep learning, each input and output variable needs to be numeric. Before fitting and evaluating a model, it has been encoded by Categorical data. One hot encoding has known as the widely accepted method, which performs optimally, if considered categorical variable does not take a huge quantity of values, in other words, you generally won't it for variables taking more than 15 different values. A new (binary) column can be generated by One hot encoding, through which the existence of each possible value from the original data can be represented.

1.7. Feature Removal

In this process, the unnecessary features (e.g. types of employment, pay scale, etc.) will get removed from the balanced dataset, as they are incapable of performing the fake news detection. Conversely, the essential features, such as title, description, requirements, department, company_profile, industry, benefits, location, required_experience, required_education, function, and fraudulent have been taken to account.

1.8. Word Embedding

Word Embedding has known being a language modeling method that helps to map the words to vectors of real numbers, through which the phrases/words in vector space can be signified alongside various dimensions. There are several approaches to execute this process, like probabilistic models,

co-occurrence matrix, neural networks, etc. The mapping of a discrete — categorical — variable to a vector of sequential numbers has termed embedding. Generally, embeddings are low-dimensional in the context of neural networks that learn discrete variable's sequential vector representations. Since these neural network embeddings possess the ability to diminish the dimensionality of categorical variables and to deliberately indicate the categories in the transformed space, they have been considered to be valuable.

1.9. Fake News Detection via Bi-Directional Long Short Term Memory (Bi-LSTM)

For post-preparation of the training and test data, a machine learning framework has trained for classifying the job posts as real and fake, for which a Bi-Directional Long Short Term Memory (Bi-LSTM) has been suggested in this study. Besides, this novel Bi-LSTM approach enhances the detection process and methodically categorizes the data as real and fake job postings. According to Zhang, C., Biś, D., Liu, X., & He, Z. (2019) Long Short-Term Memory (LSTM) has known to be a gated Recurrent Neural Network (RNN). Figure 7 demonstrates the structure of an LSTM cell, whose function can be mathematically defined as:

$$i_t = \tilde{A}(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \quad (1)$$

$$f_t = \tilde{A}(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \quad (2)$$

$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (3)$$

$$o_t = \tilde{A}(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \quad (4)$$

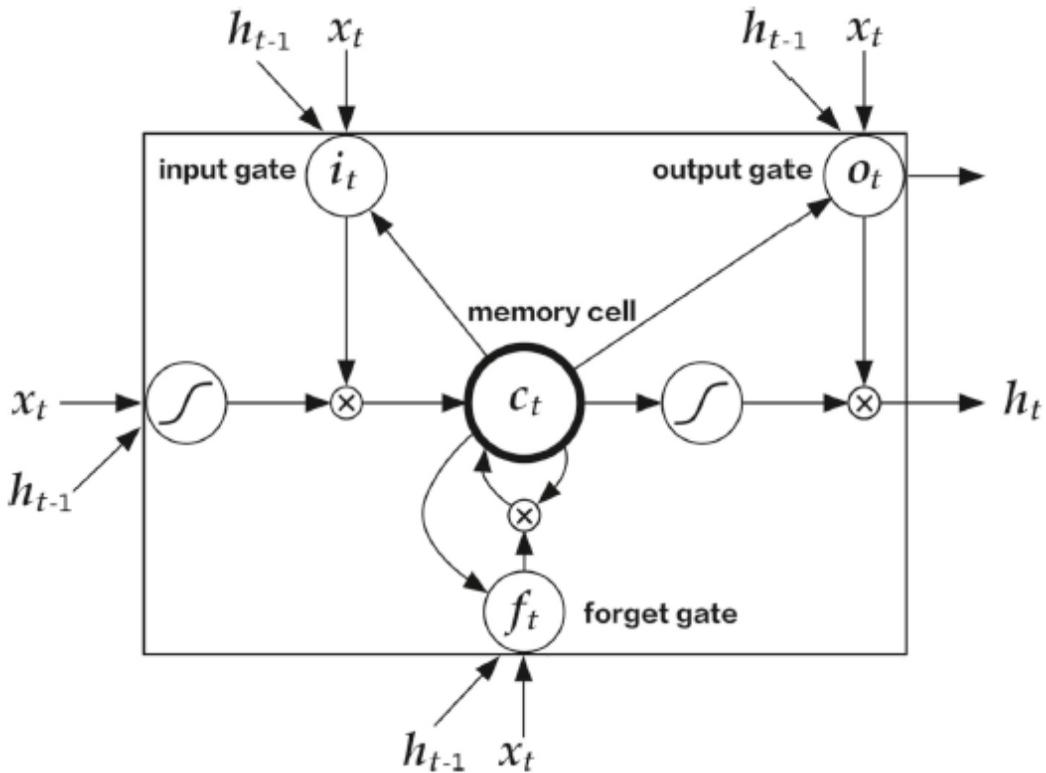
$$h_t = o_t \tanh(c_t) \quad (5)$$

Here, the Bi-LSTM classifier's weight and bias have represented by W & b .

At this point, the sigmoid function: $\tilde{A}(x) = 1 / (1 + \exp(-x))$ has indicated by A . In an LSTM cell, there are three “gates” are operating, i.e. forget gate, input gate, and output gate, which have signified by f_t , i_t , and o_t , respectively. All the three gates can be operated inside an LSTM cell alongside the trainable matrices W for retaining “valuable” info from preceding time steps and evading “invaluable” portions following the given label. This function can be a repeated task that has been presented comprehensively in prior work by Biś, D., Zhang, C., Liu, X., & He, Z. (2018).

During this study, Kiperwasser, E., & Goldberg, Y. (2016) the Bidirectional LSTM (Bi-LSTM), which has known to be such kind of LSTM has employed to detect the fake news. In the process of Bi-LSTM, the input sequence processes in both directions (forward and backward directions) accompanying independent parameters in each of the direction. In the context of multiple layers for optimal target class detection, the outputs at each time-step have concatenated, which are from each direction and turn out to be the input of the Bi-LSTM in the subsequent layer. Thereby, at each time step, the neural network node captures the overall information regarding the entire input sequence.

Figure 7. The structural form of an lstm cell according to Graves, A., & Jaitly, N. (2014)



Subsequently, the Bi-LSTM networks have involved utilizing the benefit of the aforementioned feature, through which the data correlations on both sides of the target class can get captured. As an instance, the model of the Bi-LSTM neural network can get utilized to depict the structure of the upper layer. Let the output of the first layer be $\mathbf{Y} = (y_1, \dots, y_T)$, and the output of the second layer be $\mathbf{Z} = (z_1, \dots, z_T)$, then y_i and z_i to be vectors with the same dimension D , that is to say, $y_i, z_i \in \mathbb{R}^D$. Consequently, According to Reinhart, R. F., & Steil, J. J. (2011) conclude that a two-layer Bi-LSTM with dropout delivers the optimal performance, once after enforcing different layer settings. For adjusting the output of Bi-LSTM, four optional structures have been designed to function on top of it. Consider \mathbf{H} as each structure's output, then these structures can be defined as follows,

- (i) Utilize the output directly from the Bi-LSTM, i.e. $\mathbf{H} = \mathbf{Z}$.
- (ii) Execute weighted summation within \mathbf{Y} and \mathbf{Z} . i.e. $\mathbf{H} = \lambda\mathbf{Y} + (1 - \lambda)\mathbf{Z}$, in which $\lambda \in [0, 1]$ has considered being a variable.
- (iii) Concatenate \mathbf{Y} and \mathbf{Z} along with time steps, i.e. because both \mathbf{Y} and \mathbf{Z} are $T \times D$ tensors, \mathbf{H} will be a $2T \times D$ tensor.
- (iv) Concatenate \mathbf{Y} and \mathbf{Z} along each vector \mathbf{y} and \mathbf{z} . i.e. \mathbf{H} will be a $T \times 2D$ tensor.

Post-selection of the particular upper layer structure, a max-pooling operation has implemented along time-steps over \mathbf{H} to get $\mathbf{h} \in \mathbb{R}^D$ (in case of structure (iv)). In other words, within each dimension $d \in D$ (or $2D$), the maximum value alongside the time-steps can be selected. Ultimately, adam optimizer has additionally included together with accuracy as a metric for tuning the hyperparameter.

2. RESULTS AND DISCUSSION

Based on the fake job postings dataset, the proposed Bi-LSTM classifier and prevailing classifiers, like Logistic Regression (LR), Artificial Neural Networks (ANNs), Support Vector Machine (SVM), K Nearest Neighbor (KNN) method, Random Forest (RF), and XGBoost have evaluated. In the system, the overall simulation process has been carried out in the configuration of Windows 8 Pro 64-bit OS; Intel(R) 8 GB of RAM Core(TM) i5-4260U CPU@ 1.4GHZ 2.7GHz processor. Application of algorithms has been performed through Python in Colab with Jupyter IDE. Using the Pandas Library functions from Python, the implementation of methods has carried out.

2.1. Dataset Description

From an open-access website, namely Kaggle (URL: <https://www.kaggle.com/rohan0301/predicting-fraudulency-based-on-job-posts>), the dataset has downloaded that includes 17880 job advertisements alongside human classified fake jobs, among which 17014 are real jobs, whereas 866 are fake. In that, job_id, title, location, department, salary_range, company_profile, depiction, necessities, benefits, telecommuting, has_company_logo, has_questions, employment_type, required_experience, required_education, industry, function, and fraudulent have considered being the attributes. In other words, in the first instance, the models get trained on a training set that is 80%, get tuned through validation set, eventually tested through the test set that is 20%. Python has vitally utilized for the implementation of methods.

2.2. Evaluation Metrics

Every individual model has evaluated based on various performance parameters, such as recall, precision, and f1-score. It is needed to use multiple metrics because they don't all account for the equivalent values. Precision refers to the evaluation of True Positive (TP) entities about False Positive (FP) entities. It has been formulated as,

$$Precision = \frac{TP}{TP + FP} \quad (6)$$

Recall defines the estimation of True Positive (TP) entities about False Negative (FN) entities that are uncategorized. The following Equation (7) expresses the recall estimation:

$$Recall = \frac{TP}{TP + FN} \quad (7)$$

At times, Accuracy and Recall may be inadequate to assess the performance. For example, if a mining algorithm possesses lesser recall and optimal precision, the further algorithm has necessitated for balancing the process. It is difficult to conclude that which one is optimal. For resolving this issue, F1-score has considered, which averages the recall and precision values. The estimation of F1-score can be,

$$F1 - Score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (8)$$

Accuracy refers to the number of instances that are appropriately categorised as normal and attack classes, which can be estimated as,

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}} \quad (9)$$

Macro average and the micro average have greatly utilized to compute the precision, recall, F1- Score parameters.

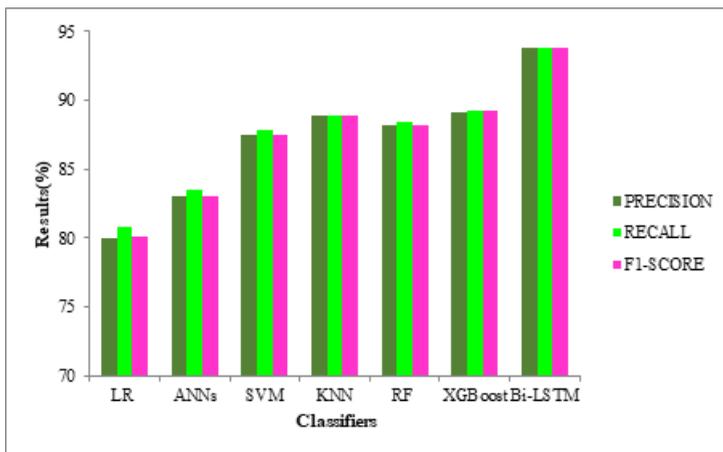
For each label with the same weight, the metrics have been estimated and summed up by Macro average.

Weighted average estimates metrics for each label and sums them up multiplied with the help of each label:

Table 2. Results comparison vs. Metrics

Average	METRICS	CLASSIFIERS						
		LR	ANNs	SVM	KNN	RF	XGBoost	Bi-LSTM
Weighted average (%)	PRECISION	80.00	82.99	87.52	88.87	88.18	89.16	93.78
	RECALL	80.79	83.50	87.76	88.87	88.39	89.23	93.78
	F1-SCORE	80.05	83.03	87.46	88.87	88.19	89.19	93.78
Macro average (%)	PRECISION	77.13	80.61	86.24	86.33	86.61	86.91	93.79
	RECALL	73.25	77.25	82.79	86.31	84.18	86.45	93.77
	F1-SCORE	74.69	78.61	84.25	86.32	85.26	86.68	93.78
ACCURACY		80.78	83.50	87.76	88.87	88.39	89.22	93.77
AUC SCORE		84.22	88.03	93.35	95.92	95.65	95.03	93.77

Figure 8. Weighted average results comparison vs. Classifiers



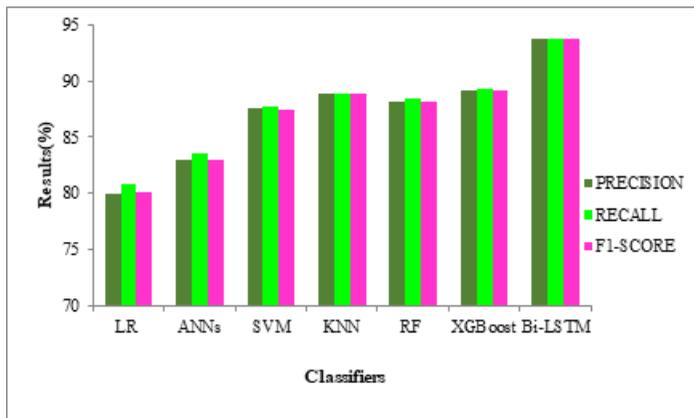
2.3. Results

In this segment, Table 2 depicts the numerical results obtained by the proposed Bi-LSTM classifier, and prevailing LR, ANNs, SVM, KNN, RF, and XGBoost classifiers, in terms of the aforementioned metrics.

In figure 8, the graphs compare the weighted average results of precision, recall, and F1-score obtained by the proposed Bi-LSTM classifier and prevailing LR, ANNs, SVM, KNN, RF, and XGBoost classifiers. Among the conventional machine learning classifiers, the Bi-LSTM classifier proves to be efficient to provide optimal performance. Accordingly, it secures a 93.78% weighted average F1-score, which is superior to other prevailing approaches, since 80.05%, 83.03%, 87.46%, 88.87%, 88.19%, and 89.19% of F1-score have attained by LR, ANNs, SVM, KNN, RF, and XGBoost, respectively (as represented in Table 2). It can be observed from the graph that if the length of the dataset increases, the Bi-LSTM classifier also improves linearly from the job postings dataset. Thus, it can get concluded that the proposed Bi-LSTM classifier is a promising method for fake detection, as it delivers efficient performance on metrics that assures claiming attention for investigating a bigger dataset, in the future.

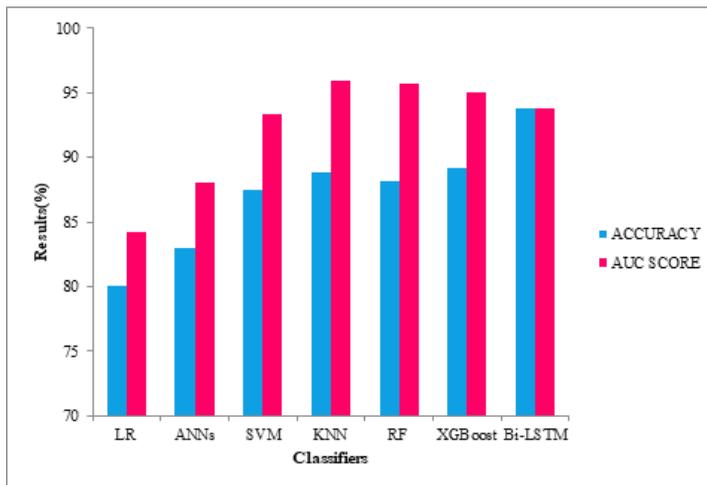
Figure 9 compares the macro-average outcomes of precision, recall, and F1-score obtained by the suggested Bi-LSTM classifier and prevailing LR, ANNs, SVM, KNN, RF, and XGBoost classifiers. From the graphs, it can be observed that the proposed Bi-LSTM classifier has the efficiency to procure a 93.78% macro average F1-score, which is comparatively higher than LR, ANNs, SVM, KNN, RF, and XGBoost methods, as they solely attain 74.69%, 78.61%, 84.25%, 86.32%, 85.26%, and 86.68% of F1-score, correspondingly. Eventually, with a larger dataset, the Bi-LSTM classifier proves to be proficient for delivering the optimal performance in terms of macro-average results.

Figure 9. Macro average results comparison vs. Classifiers



In Figure 10, the performance of the proposed Bi-LSTM classifier has compared with prevailing LR, ANNs, SVM, KNN, RF, and XGBoost classifiers in terms of accuracy and AUC. During that, the proposed classifier proves its efficiency to obtain 93.77% accuracy, which is considerably superior to the existing LR, ANNs, SVM, KNN, RF, and XGBoost methods, since they solely attain 80.78%, 83.50%, 87.76%, 88.87%, 88.39%, and 89.22%, respectively. Ultimately, with a batch size of 32 and utilizing two epochs, the model runs for attaining the optimal accuracy of 94%.

Figure 10. Accuracy and auc results comparison vs. Classifiers



3. CONCLUSION AND FUTURE WORK DIRECTION

This study proposes the classifier called Bi-Directional Long Short Term Memory (Bi-LSTM) for the initial stage detection of job-associated fake news that has been posted over the social media. It includes various features, specifically the company profile from the job description to classify the job posted on social media as either fake or real. There are two significant phases involved in the proposed framework, i.e. Natural Language Processing (NLP) and Text conversion and Classification. Among that, the first phase has further classified into many pre-processing steps, like missing values, visual representation of text data (Word Cloud), fill null values, cleaning of text features through stop words removal, tokenization. Besides, lemmatizing has executed through wordnet as regards the cleaning of a text document. Consequently, conversion from text data to numerical format has carried out using one-hot encoding, besides word embedding has also been presented to map the numbers, such as words to vectors and real numbers. During the third phase, the proposed Bi-LSTM classifier has presented to classify the fake job posting, where the input sequence processes in both of the directions such as forward and backward direction by accompanying independent parameters in each direction. In the instance of multiple layers for optimal target class detection, the outputs at each time-step have concatenated, which are from each direction and turn out to be the input of the Bi-LSTM in the subsequent layer. Even though the company profile has been vigorously recommended by Bi-LSTM classifier as a robust feature, yet it additionally focuses on some other valuable factors, such as requirements, benefits, and description of the job, since they significantly possess the predictive ability. Through this, the capability of the proposed method to obtain maximum rates of accuracy, precision, recall, F1-score has been depicted, which represents that the Bi-LSTM is highly efficient in the detection of the fake job advertisement. Alongside basic classifiers, the proposed Bi-LSTM provides highly optimal outcomes accompanied by hyper-parameter tuning.

To lessen the errors from each specific variable, in place of the single classifier, an ensemble model has been employed to make it a robust framework. In the future, this study can get further stretched by focusing on the experiment over a bigger dataset vs lower computational time in the fake news detection. During the process, a few Artificial Intelligence (AI) techniques, specifically Harmony Search (HS), Ant Colony Optimization (ACO), Swallow Swarm Optimization (SSO), Simulated Annealing (SA), Big Bang–Big Crunch (BB–BC), and Particle Swarm Optimization (PSO) can be considered to implement the feature removal.

REFERENCES

- Alghamdi, B., & Alharby, F. (2019). An intelligent model for online recruitment fraud detection. *Journal of Information Security*, 10(03), 155–176. doi:10.4236/jis.2019.103009
- Ibrishimova, M. D., & Li, K. F. (2019). A machine learning approach to fake news detection using knowledge verification and natural language processing. *International Conference on Intelligent Networking and Collaborative Systems*, 223-234.
- Zhou, X., Zafarani, R., Shu, K., & Liu, H. (2019). Fake news: Fundamental theories, detection strategies and challenges. *Proceedings of the twelfth ACM international conference on web search and data mining*, 836-837.
- Khan, J. Y., Khondaker, M., Islam, T., Iqbal, A., & Afroz, S. (2019). *A benchmark study on machine learning methods for fake news detection*. arXiv preprint arXiv:1905.04749.
- Ahmed, H., Traore, I., & Saad, S. (2017). Detection of online fake news using n-gram analysis and machine learning techniques. *International conference on intelligent, secure, and dependable systems in distributed and cloud environments*, 127-138. doi:10.1007/978-3-319-69155-8_9
- Kaliyar, R. K. (2018). Fake news detection using a deep neural network. *International Conference on Computing Communication and Automation (ICCCA)*, 1-7. doi:10.1109/CCAA.2018.8777343
- Singhania, S., Fernandez, N., & Rao, S. (2017, November). 3han: A deep neural network for fake news detection. *International Conference on Neural Information Processing*, 572-581. doi:10.1007/978-3-319-70096-0_59
- Chen, K., Wang, J., Chen, L. C., Gao, H., Xu, W., & Nevatia, R. (2015). *Abc-cnn: An attention based convolutional neural network for visual question answering*. arXiv preprint arXiv:1511.05960.
- Agarwal, V., Sultana, H. P., Malhotra, S., & Sarkar, A. (2019). Analysis of Classifiers for Fake News Detection. *Procedia Computer Science*, 165, 377–383. doi:10.1016/j.procs.2020.01.035
- Poddar, K., & Umadevi, K. S. (2019). Comparison of various machine learning models for accurate detection of fake news. *Innovations in Power and Advanced Computing Technologies (i-PACT)*, 1, 1-5.
- Mahabub, A. (2020). A robust technique of fake news detection using Ensemble Voting Classifier and comparison with other classifiers. *SN Applied Sciences*, 2(4), 1–9. doi:10.1007/s42452-020-2326-y
- Thota, A., Tilak, P., Ahluwalia, S., & Lohia, N. (2018). Fake news detection: A deep learning approach. *SMU Data Science Review*, 1(3), 1–21.
- Zhang, J., Dong, B., & Philip, S. Y. (2020). Fakedetector: Effective fake news detection with deep diffusive neural network. *IEEE 36th International Conference on Data Engineering (ICDE)*, 1826-1829.
- Yang, Y., Zheng, L., Zhang, J., Cui, Q., Li, Z., & Yu, P. S. (2018). *TI-CNN: Convolutional neural networks for fake news detection*. arXiv preprint arXiv:1806.00749.
- Monti, F., Frasca, F., Eynard, D., Mannion, D., & Bronstein, M. M. (2019). *Fake news detection on social media using geometric deep learning*. arXiv preprint arXiv:1902.06673.
- Zhang, J., Cui, L., Fu, Y., & Gouza, F. B. (2018). *Fake news detection with deep diffusive network model*. arXiv preprint arXiv:1805.08751.
- Liu, Y., & Wu, Y. F. (2018). Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 354–361.
- Roy, A., Basak, K., Ekbal, A., & Bhattacharyya, P. (2018). *A deep ensemble framework for fake news detection and classification*. arXiv preprint arXiv:1811.04670.
- Dutta, S., & Bandyopadhyay, S. K. (2020). Fake Job Recruitment Detection Using Machine Learning Approach. *International Journal of Engineering Trends and Technology*, 68(4), 48–53. doi:10.14445/22315381/IJETT-V68I4P209S
- Millstein, F. (2020). *Natural language processing with python: natural language processing using NLTK*. Frank Millstein.

Hardeniya, N., Perkins, J., Chopra, D., Joshi, N., & Mathur, I. (2016). *Natural language processing: Python and NLTK*. Packt Publishing Ltd.

Bird, S., Klein, E., & Loper, E. (2009). *Natural language processing with Python: analyzing text with the natural language toolkit*. O'Reilly Media, Inc.

Perkins, J. (2010). *Python text processing with NLTK 2.0 cookbook*. Packt Publishing Ltd.

Zhang, C., Biś, D., Liu, X., & He, Z. (2019). Biomedical word sense disambiguation with bidirectional long short-term memory and attention-based neural networks. *BMC Bioinformatics*, 20(16), 1–15. doi:10.1186/s12859-019-3079-8 PMID:31787096

Biś, D., Zhang, C., Liu, X., & He, Z. (2018). Layered multistep bidirectional long short-term memory networks for biomedical word sense disambiguation. *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 313-320. doi:10.1109/BIBM.2018.8621383

Kiperwasser, E., & Goldberg, Y. (2016). Simple and accurate dependency parsing using bidirectional LSTM feature representations. *Transactions of the Association for Computational Linguistics*, 4, 313–327. doi:10.1162/tacl_a_00101

Reinhart, R. F., & Steil, J. J. (2011). A constrained regularization approach for input-driven recurrent neural networks. *Differential Equations and Dynamical Systems*, 19(1-2), 27–46. doi:10.1007/s12591-010-0067-x

Graves, A., & Jaitly, N. (2014). Towards end-to-end speech recognition with recurrent neural networks. *PMLR International Conference on Machine Learning*, 1764-1772.