Afaan Oromo Fake News Detection Using Natural Language Processing and Passive-Aggressive

Daraje kaba Gurmessa

Mizan-Tepi University, Ethiopia

Email: darajekaba2020@gmail.com and darajekaba@mtu.edu.et

Abstract— The main objective of this study is to develop Afaan Oromo fake news detection system. The designed system involves preprocessing like tokenization, Normalization, stop word removing and abbreviation resolving, feature extraction-like Term-Frequency-inverted document frequency, term frequency, and hash to know word importance that appears in the news and word appears in the corpus and N-grams which are a powerful Natural Language Processing technique in order to capture semantic and syntactic sequences was also used. All possible combination of features extraction techniques and natural processing techniques were used with a passiveaggressive classification algorithm. Passive-Aggressive performs 97.2% with a classification error of 2.8% which was better than ensemble algorithms like gradient boosting and random forest and linear classifier like multinomial Naïve Bayes. Finally, a python Django was used for the web-based deployment of the model system using the Term Frequency-Inverted Document Frequency feature extraction with unigram and Passive aggressive classification algorithm.

Keywords— Oromo, Fake News Detection, Passive-Aggressive, NLP, TF-IDF.

1. INTRODUCTION

Ethiopia is the third-ranked country in the world in internet freedom based on the 2018 net freedom score (Figure 1) and is the first-ranked horn African country with the largest nation nationalities and peoples having different cultures, languages, heritage and political view [1]. So fake news can be easily subject to the spark of violence like race, culture, religion, and language by which peoples identify themselves. The population of Oromo is more than 40 million in Ethiopia and 3rd largest single nationality group in Africa [2]. The Oromo nation has a single common mother tongue, called the Oromo language or Afaan Oromo. It is the third most widely spoken language in Africa as a mother tongue, next to Hausa and Arabic. Today, Afaan Oromo is serving as an official language of Oromia regional state which is the largest regional state among the current federal states of Ethiopia. Being an official language, it is also used as a medium of instruction for primary and junior secondary schools of the region [3]. It is also a field of specialization at Diploma, Bachelor Degree, and Master's Degree levels at various universities in Oromia regional state.

Fake news is a series problem in international communities due to its negative impact on society and individuals. Fake news is not a new phenomenon. It is also a problem during the first and second World Wars [4]. Nowadays it is acting as a war weapon all over the world due to the platform of social media. Fake news is different from spam, rumors, and yellow journalism. It is "a completely fabricated claim or story with an intention to deceive, often for secondary gain" [5]. The secondary gain in a sense fake news exists for both political advantage and financial generation or fundraising; different media write the interest of rich persons, but this may be not the only reason. So; the reasons behind fake news include media manipulation and propaganda, political and social influence, provocation, social unrest, and financial profit.



Figure: 1 Freedom on the Net Improvements and Declines 2018 (source: FOTN_2018)

2. RELATED WORKS

The paper [6] was discussed application of natural language processing techniques with multinomial naïve Bayes for the detection of "fake news" on 752 news

datasets that prepared for Afaan Oromo language. They got term frequency of unigram of their model identifies fake news with an accuracy of 96%. The amount of dataset [6] used for this study is relatively larger than the news dataset prepared for the Indonesian language [7] and Russian language [8] but smaller when compared to the dataset for English so it requires further development.

3. METHODOLOGY

3.1 System architecture

Afaan Oromo fake News Detection System has three major components like many other fake news Detection System. These are Search Engine Component, Preprocessing Component and Detection Component. However, these components have subcomponents that are unique for the language they are designed for.



Figure: 2 Afaan Oromo Fake News Detection System Architecture [6]

3.2 Data Exploration

In this section how to take a quick look at the data and get a feel for its contents is done. To do so, the researcher uses a Pandas Data Frame and check the shape, head and apply any necessary transformations. For every news, the following information is available; Article text, Article address, Article title and Article label (fake or truthful). In This news dataset 571,738,999,696 of terms are there, one news has a maximum of 10704 terms, a minimum of 28.000000 terms, 760.291223 mean, and 975.911 standard deviation.

Once collected, prepared and preprocessed in a required format a number of experiments were conducted for comparing models accuracy for the proper choice in developing system prototype. The word cloud in (Figure 3 and Figure 4) shows which word is a more frequent word in the whole dataset before removing stop words and after removing stop words respectively. (Figure 3) shows that the frequent words ('fi', 'akka','kan' and etc.) are seen bold because it is before removing the stop words from the corpus. (Figure 4) shows the cloud of the same data in (Figure 3) but the stop words are removed. It shows words like ('hojii', 'nama', 'motummaa, 'tokko' and etc.) are seen as bold. Additionally, the frequency of words in fake labeled news datasets and real labeled news dataset is different. For example, the word 'hin' is negation which has 'not' English equivalent meaning is more common in fake news dataset (Figure 6) news than real news datasets (Figure 5).

The wordcloud of the complete all dataset before stopword removed



Figure:3 The Word Cloud of the All Dataset Before Removing Stop Words The wordcloud of the complete all dataset After stopword removed



Figure 4 The Wordcloud of All Dataset After Removing Stopwords



Figure:5 The Wordcloud of the Real Labeled Dataset of All Dataset



Figure 6: The Wordcloud of the Fake Labeled Dataset of All Dataset

3.3 Feature extraction

There are several ways of feature extraction [9]. Different approaches on the same dataset are compared to determine which method gives the best accuracy was experimented. Three features extraction methods are applied in this research for comparing, namely: - Hashvectorizer, countyectorizer and Term Frequency-Inverted Document Frequency (TF-IDFvectorizer).

3.2.1. Term Frequency-Inverse Document Frequency (TF-IDF)

Term Frequency-Inverse Document Frequency is also a method used to represent text in a format that can be easily processed by the machine learning algorithms. It is a numerical statistic that shows how important a word is to news in a news dataset. The importance of a word is proportional to the number of times the word appears in the news (fake and real) but inversely proportional to the number of times the word appears in the news dataset (fake or real) [10].

In this Research, a computer learns, how to read and understand the differences between real news and fake news using Natural Language Processing (NLP). This is done by using TF-IDFvectorizer, countVectorizer, and hashVectorizer. TF-IDF is used to determine word importance in a given article in the entire news dataset. The frequency of the words is rescaled by considering how frequently the words occur in all the news dataset. Due to this, the scores for frequent words are also frequent among all the documents are reduced. This way of scoring is known as Term Frequency – Inverse Document Frequency. A term importance increases with the number of times a word appears in the document, however, this is counteracted by the frequency of the word in the corpus.

Let C denote a corpus, or set of news documents. Let n denote a news document $n\in$; a news document is defined as a set of words w. Let $m_w(d)$ denote the number of times word w appears in news document d. Hence, the size of news d is

$$|n| = \sum\nolimits_{w \in n} m_w(n)$$

The normalized Term Frequency (TF) for word w with respect to news document d is [10] defined as follows:

$$TF(W)_n = \frac{m_n(w)}{|n|}$$

The Inverse Document Frequency (IDF) for a term w with respect to news corpus C, denoted $IDF(W)_C$ is the logarithm of the total number of news in the corpus divided by the number of news where this particular term appears and computed as follows: [10]

$$IDF(W)_{\mathcal{C}} = 1 + \log(\frac{|\mathcal{C}|}{|\{n: \mathcal{C}|w \in n\}|})$$

One of the main characteristics of IDF is it weights down the term frequency while scaling up the rare ones. For example, words such as "jiru" and "ture" often appear in the text, and if we only use TF, terms such as these will dominate the frequency count. However, using IDF scales down the impact of these terms.

TF-IDF for the word w with respect to news d and corpus D is calculated as follows:

 $TF - IDF(W)_{n,C} = TF(W)_n \times IDF(W)_C$

So for example, let say we have a piece of news with 800 words and we need the TF-IDF for the word "hidhan". Assuming that the word "hidhan" occurs in the document 16 times then $TF = \frac{16}{800} = 0.02$ then the IDF could be calculated; let's assume that there are 752 documents and "hidhan" appears in 188 of them then IDF (hidhan) = $1 + \log(\frac{752}{188}) = 1.6$. Then TF - IDF (hidhan) = $0.2 \times 1.6 = 0.32$

3.2.2. Classification Reports:

From the above confusion matrix described in FIGURE 7 the classification report was calculate

PAC-TFIDF

Table 3: Classification Report for PAC-TFIDF

Classification Report	Precision	Recall	F1Score	Support
Fake	95%	97.9%	96%	96
Real	98.6%	96.6%	97%	146

3.4 Passive Aggressive Classifier

Passive Aggressive algorithm is a margin-based online learning algorithm for binary classification. It is also an algorithm of a soft margin-based method and robust to noise. [11]

Let's suppose to have a dataset:

 $X = \{\overline{x_0}, \overline{x_1}, \dots, \overline{x_t}\} \text{ where } \overline{x_i} \in \mathbb{R}^n$

 $Y = \{\overline{y_0}, \overline{y_1}, \dots, \overline{y_t}\}$ where $y_i \in \{-1, +1\}$

The index t has been chosen to mark the temporal dimension. In this case, in fact, the samples can continue arriving for an indefinite time. Of course, if they are drawn from the same data generating distribution, the algorithm will keep learning (probably without large

parameter modifications), but if they are drawn from a completely different distribution, the weights will slowly *forget* the previous one and learn the new distribution. For simplicity, let us also assume it is working with a binary classification based on bipolar labels. A Passive-Aggressive algorithm works generically with this update rule: [11]

$$\begin{cases} \overline{w}_{t+1} = argmin_{\overline{w}} \frac{1}{2} \|\overline{w} - \overline{w}_t\|^2 + c\xi^2 \\ L(\overline{w}; \overline{x_t}, y_t) \le \xi \end{cases}$$

Let's assume the slack variable $\xi = 0$ (and *L* constrained to be 0). If a sample *x* (*t*) is presented, the classifier uses the current weight vector to determine the sign. If the sign is correct, the loss function is 0 and the *argmin* is *w* (*t*).

This means that the algorithm is passive when a correct classification occurs. Let's now assume that a misclassification occurred: The angle $\theta > 90^\circ$, therefore, the dot product is negative and the sample is classified as -1, however, its label is +1.

In this case, the update rule becomes very aggressive, because it looks for a new w which must be as close as possible as the previous (otherwise the existing knowledge is immediately lost), but it must satisfy L = 0 (in other words, the classification must be correct).

The introduction of the slack variable allows one to have soft-margins (like in SVM) and a degree of tolerance controlled by the parameter *C*. In particular, the loss function has to be $L \leq \xi$, allowing a larger error. Higher *C* values yield stronger aggressiveness (with a consequent higher risk of destabilization in the presence of noise), while lower values allow a better adaptation. In fact, this kind of algorithms, when working online, must cope with the presence of noisy samples (with wrong labels).

A good robustness is necessary, otherwise, too rapid changes produce consequent higher misclassification rates. After solving both update conditions, it gets the closed-form update rule: [11]

$$\overline{w}_{t+1} = \overline{w}_t + \frac{\max(0, 1 - y_t(\overline{w}^T \cdot \overline{x}_t))}{\|x_t\|^2 + \frac{1}{2c}} y_t \overline{x}_t$$

This rule confirms the expectations: the weight vector is updated with a factor whose sign is determined by y(t)and whose magnitude is proportional to the error. Note that if there's no misclassification the nominator becomes 0, so w(t + 1) = w(t), while, in case of misclassification, w will rotate towards x(t) and stops with a loss $L \le \xi$.

This classification algorithm is used in combination with different feature extraction algorithms. It was used as a combination of the Passive-Aggressive classifier with TFIDF at the n-gram level. As mentioned (above) in section (3.4), machine learning classification algorithms were studied. The algorithms were used to create learning models, and then, the learned models were used to predict the labels assigned to the testing data. Experiment results were then presented, analyzed, and interpreted on section (4.1, 4.2, 3.2.1, 4.3, and 4.1) below.

4.1. Experiment 1: Accuracy of Models Using 752 Datasets

By increasing dataset size to compare each previous possible combination of classifier, feature selections, and n-grams. So using a 33% news test set among 752 news datasets the best accuracy score of 97.2% was obtained using passive-aggressive, TF-IDF and unigram.

Features	Classifier	N-gram	Score
selection			
Term Frequency	Multinomial	Unigram	96
Term Frequency	Multinomial	Bigram	92.8
Term Frequency	Passive-aggressive	Unigram	95.6
Term Frequency	Passive-aggressive	Bigram	92.8
Term Frequency	Gradient Boosting	Unigram	95.2
Term Frequency	Gradient Boosting	Bigram	81.5
Term Frequency	Random Forest	Unigram	96.4
Term Frequency	Random Forest	Bigram	83.9
TF-IDF	Multinomial	Unigram	95.6
TF-IDF	Multinomial	Bigram	93.2
TF-IDF	Passive-aggressive	Unigram	97.2
TF-IDF	Passive-aggressive	Bigram	94
TF-IDF	Gradient Boosting	Unigram	79.1
TF-IDF	Gradient Boosting	Bigram	79.9
TF-IDF	Random Forest	Unigram	93.6
TF-IDF	Random Forest	Bigram	91.6
Hash	Multinomial	Unigram	83.5
Hash	Multinomial	Bigram	92.4
Hash	Passive-aggressive	Unigram	96
Hash	Passive-aggressive	Bigram	91.6
Hash	Gradient Boosting	Unigram	96.4
Hash	Gradient Boosting	Bigram	82.3
Hash	Random Forest	Unigram	92.4
Hash	Random Forest	Bigram	86.3
	Features selection Term Frequency Term Frequency Term Frequency Term Frequency Term Frequency Term Frequency Term Frequency Term Frequency TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF TF-IDF Hash Hash Hash Hash Hash Hash Hash Hash	Features Classifier selection	Features Classifier N-gram selection Nultinomial N-gram Term Frequency Multinomial Bigram Term Frequency Passive-aggressive Unigram Term Frequency Passive-aggressive Bigram Term Frequency Gradient Boosting Unigram Term Frequency Gradient Boosting Bigram Term Frequency Gradient Boosting Bigram Term Frequency Random Forest Unigram Tern Frequency Random Forest Unigram TF-IDF Multinomial Unigram TF-IDF Multinomial Bigram TF-IDF Multinomial Bigram TF-IDF Gradient Boosting Unigram TF-IDF Random Forest Unigram TF-IDF Random Forest Unigram TF-IDF Random Forest Bigram Hash Multinomial Bigram Hash Passive-aggressive Unigram Hash Gradient Boosting Bigram Hash Gradient Boosting Bigram Hash Gradient Boosting Unigram Hash Gradient Boosting Unigram Hash Gradient Boosting

Table: 1 Classification Accuracy for Experiment 2

Classification Accuracy at large data set size is better than the small data set size and at Word level performed better than bi-gram level as we can see from the above (aggressive, TF-IDF and unigram. Table (1).

The accuracy for Gradient Boosting with TF-IDF at word level was the lowest at 79.1% while Linear Passive Aggressive Classifier, using TF-IDF vectors at word level performed 97.2 compared to other classifier but when it is combined with bigram in place of unigram the accuracy was lowered to 94%. So using unigram shows the best performance than bigram.

Since Classification accuracy alone is not enough to determine the effectiveness of the model; other metrics were also explored especially for two algorithms (Multinomial and Passive-Aggressive) at the word level, using TF-IDF Vectors.

In another experiment, the parameters described above are included. With the increase of MDM from 0 to 1 in the step of 0.1, the classification accuracy of the two models increased significantly as depicted by table below (Table 2).

Classification Accuracy		Multinomial	Passive-Aggressive Classifier				
Maximum Document Frequency	0.1	95.9	93.5				
	0.2	96.7	93.5				
	0.3	96.7	93.5				
	0.4	95.5	94				
	0.5	95.5	94				
	0.6	95.5	95				
	0.7	96.7	97.2				
	0.8	95.5	97				
	0.9	95.5	93.5				

The best performing model was Multinomial with 96.7 % at MDM(X = 0.7) and Passive-Aggressive Classifier with 97.2%. Beyond, 0.7 the algorithms did not show improvement. So, MDM with 0.7 was chosen as the optimal value for Multinomial and also for passive-aggressive. Henceforth, the obtained the Classification reports including precision, recall, fscore of both models at is calculated using MDM(X=0.7)

4.2. Confusion Matrix

A confusion matrix shows the proper labels on the main diagonal (top left to bottom right). The other cells show the incorrect labels, often referred to as true negative or false negatives. Depending on the problem, one of these might be more significant. For example, for the fake news problem, it is more important than they don't label real news articles as fake news if so, the researcher might want to eventually weight the accuracy score to better reflect this concern. This was done across three random tests and thirty-three percent test set of the



Figure 7: Passive-Aggressive with TF-IDF Vector and Unigram, Bigram Respectively



Figure 8: Multinomial with TF-IDF Vector and Unigram, Bigram Respectively



Figure 9: Random Forest, TF-IDF Vectorizer with Unigram and Bigram Respectively

Sensitivity is how sensitive the classifier is to detect fake news, while Specificity is how selective or specific the model is in predicting real news.

Choosing the metric depends on what kind of application is going to be developed. The positive class in this binary classification is class "Fake". Therefore, Sensitivity should be higher, because false positives are more acceptable than false negatives in classification problems of such applications. The sensitivity is high for both the models and is having equal value. By optimizing more for Sensitivity, we can get better results. By decreasing the threshold to some extent for predicting fake news, the Sensitivity of the classifier can increase. This would increase the number of True Positives. In this work, a threshold is set to 0.7 by default but the researcher can adjust it to increase sensitivity or specificity depending on what is wanted to be.

Precision value for Linear PAC-TFIDF at 95% is higher than Multinomial -TFIDF, which is 91% and Recall values (Sensitivity) of PAC-TFIDF was calculated as 97.9% which is higher than Multinomial -TFIDF, which is 97.8% for both models. F1 score for PAC-TFIDF is 96% and the f1 score for Multinomial –TFIDF is 94%. Classification Error: It means overall, how often the model is incorrect, also called as Misclassification Rate. Classification Error for Linear PAC-TFIDF = 100–97.2 = 2.8%

4.3. Receiver Operating Characteristics Curve

It is a way to check how various thresholds affect sensitivity and specificity, without actually changing the threshold.



Figure 10: Roc Curve for Passive Aggressive and Multinominal with TFIDF Using Unigram and Bigram Respectively

4.4. Other Experiments

In this experiment, in the frequency of words signifying Fakeness was calculated in (Figure 11). The general idea is that if there is a number of such words used in news then that news has a high probability of being Fake. If no such words are present, then that article is most probably a real article.

•	•	
FAKE	-1.4763011952250658 wbo	REAL 1.0267866216333714 sirna
FAKE	-1.4028038991067868 waraanni	REAL 0.9494083343861881 beeksise
FAKE	-1.355984126477589 osoo	REAL 0.9486871936794211 keessatti
FAKE	-1.2916476888165664 guyyaa	REAL 0.9228585171079967 itiyoophiyaa
FAKE	-1.167678174368165 malee	REAL 0.9013682704760803 2011
FAKE	-1.1491388502232995 vemen	REAL 0.898219788959887 jedhan
FAKE	-1.0720437391579463 hardhaa	REAL 0.8941648096740794 gara
FAKE	-1.0505007407250362 wbon	REAL 0.8390956156634011 finfinnee
FAKE	-1.0495847183550668 zoonii	
EAVE	4 0060040000464604 0000000	REAL 0.83134/0684390533 Wallin
FAKE	-1.006224930101601 deerroon	REAL 0.8208567696454153 adda
		REAL 0.809439992724103 addunyaa
FAKE	-0.9841281160896643 lolli	REAL 0.8067063979476597 turan
FAKE	-0.9509634752761226 oromoof	REAL 0.8030719658322774 bara

Figure 11: Top Informative Tokens for Passive

Aggressive

As clear from the above (Figure 11), the "fake" word has the maximum negative TF-IDF value of -1.47 in the Fake class and maximum positive TF-IDF value of 1.026 in the Real class. So this result shows the distance between Fake and Real claims.

4.5. Discussion

This section discusses the results of various experiments. The type of data used for training and the size of features that affect the classifier performance. It was able to observe an improvement in the accuracy results in experiment 1 (above), which was a bigger dataset with 752 news dataset items and it is bigger than the experiment which is only with 200 news datasets.

As observed from the result an increase in the n-gram size would cause a decrease in accuracy. In both news datasets one and two, the unigram and bi-gram was performed better than the trigram and quad-gram.

Resolving abbreviations, removing words that exist in both fake and real labeled news more than 0.7 (MDM) and counting synonyms words as one-word increase the performance of the system. But, stemming the words decrease the performance of the system because of the meaning of the word completely changed by adding affixes and removing the affixes in Afaan Oromo.

The semantic measurement achieved rather good results. The researcher was able to identify near-duplicated content with up to 55% of the text changed. In reality, the majority of near-duplicated contents are less than 35% different from the original contents. The faker tends to report their early work and only change certain words most of the time. Thus, it was believed that the semantic measurement approaches able to classify nearduplicate correctly in real-world data.

Generally, from the above experiment results (4.1, 4.2, 3.2.1, 4.3 and 4.1), was concluded that passive-aggressive algorithm, using Term-Frequency Inverse Document Frequency vector (Word level) at maximum

document frequency of 0.7, gave the best performance. Based on the above experiment results; finally, it was chosen as the best model to determine the truth of Afaan Oromo news in social media. So based on this conclusion the system prototype was developed and the user interface screenshot sample is shown in (Figure 12).



Figure 12: System Interface for Afaan Oromo Fake News Detection.

4.6. Conclusion

The purpose of the study is to develop a Content-Based Afaan Oromo fake news detection system using a machine learning approach that can serve as a basic building block for Fake news detection. Therefore, by exploring more social media features in the experiments, and by combining them the researcher created an effective and reliable system for detecting Fake news. To-do so data was collected manually from Facebook pages and accounts by journalism experts and labeled as fake, real and unverified. Only fake and real labeled dataset was selected as a training and testing set because they clearly distinct fake from real. The designed system for Afaan Oromo fake news detection system involves preprocessing like tokenization, Normalization, stop word removing and abbreviation resolving, feature extraction using like Term-Frequency-inverted document frequency, term frequency, and hash to know word importance that appears in the news and word appears in the corpus and N-grams which are a powerful Natural Language Processing technique in order to capture semantic and syntactic sequences was also used. Based on features extracted different classification algorisms like multinomial Naive Bayes, random forest, gradient boosting, and passive-aggressive are used. The final model was created by combining the n-grams, features extracted and transformed, and classifiers. The performance of the models was accessed and compared on the same news dataset using the most significant metrics by which a machine learning model performance is measured like classification accuracy, Error matrix, Classification Report (precision and recall) and area under Receiver Operating Characteristics Curve.

Since the dataset is a great issue in the Afaan Oromo language, the classification was tested on a small number of news dataset items. As was shown in experiment (4.1) above adding more data to the news dataset test the consistency of the performance thereby increasing the trust of users on the system. More linguistic-based features were applied to the dataset to determine the news truth confidence score. Even though the dataset is a great issue the model Linear Passive-Aggressive with Term Frequency-Inverted Document Frequency vector and unigram performs unexpected with the highest accuracy of 97.2%, sensitivity of 97.9% and ROC AUC score of 97.5%.

The model served as a better model as compared with others listed on aggressive, TF-IDF and unigram.

Table 1. Because the passive-aggressive algorithm is an online algorithm and fake news detection is also an online challenging problem both fits each other.

The model generated some errors. Indeed, it is possible to anticipate such considerable contributions and positive effects of the system since Afaan Oromo is one of the morphologically rich and complex languages. The error rate was about 2.8%.

This shows that the system can be performed with low error rates in high inflected languages such as Afaan Oromo.

4.7. Contributions of the Study

The main contributions of this study are summarized as follows: -

- 1. To develop text content-based Fake news detection system a number of corpus is vital. Accordingly, the study prepared labeled Afaan Oromo news corpus.
- 2. The general architecture of text content-based Afaan Oromo fake news detection is proposed.
- 3. The study identifies basic challenges in developing text content-based fake news detection systems and the possible strategies to solve those challenges.
- 4. Comparing supervised machine learning approaches by considering feature extraction methods and language features.
- 5. As fake news is becoming a serious problem in social media, this study paves a way for developing fake news detection in Afaan Oromo for user confidence.
- 6. The research critically reviewed the performance evaluation measurement for AFND systems and identified that there is no standard benchmark and common measure used to evaluate the system. Mostly they are based on the dataset used and type of classification algorisms used.

4.8. Future Work

1. Based on experimental results and analysis the researcher recommends the following observed

points to be taken into consideration for future work to enhance the effectiveness of Afaan Oromo content-based fake news detection system. The researcher would like to recommend: -

- 2. Incorporating in a future model, statistical features, and features that reflect the writer's styles such as the number of slang words or filler words in the text.
- 3. Finding a standard corpus and gathering real news that almost appears as Fake news will improve the training of the model.
- 4. Articles are a primary source of information. There is a massive amount of data in the real world. So testing models on big data, to investigate how it will perform. Furthermore, given that high dimensionality issues will arise so it needs to explore the effect of feature selection techniques such as information gains and chi-square, in our future work.
- 5. The next researcher can also apply a recurrent neural network by combining a content-based or linguistic-based and social context to detect fake news on social media.
- The next researcher can use shallow and deep syntax analysis which uses POS (parts-of-speech) tags and Probabilistic Context-Free Grammars (PCFG) respectively and named entity recognition.

5. REFERENCES

- [1] Worldatlas, "Most Ethnically Diverse Countries in the World," 18 june 2019. [Online]. Available: https://googleweblight.com/i?u=https://www.worl datlas.com/articles/most-ethnically-diversecountries-in-the-world.html&hl=en-ET.
- [2] CenteralStatisticalAgency, "2007 population and housing census of Ethiopia," Federal Democratic Repubilic of Ethiopia, Addis Abeba, 2012.
- [3] Girma, "Afaan Oromo news text summarizer," Addis Ababa University, Addis Ababa, 2012.
- [4] Sakha, "https://www.sakhaglobal.com," 03 05 2019. [Online]. Available: https://www.sakhaglobal.com/index.php/2018/09/ 28/detecting-fake-news-through-nlp/.
- [5] FakeNewsChallenge, "Exploring how artificial intelligence technologies could be leveraged to combat fake news.," 2019. [Online]. Available: http://www.fakenewschallenge.org/. [Accessed 01 10 2019].
- [6] K. Daraje, M. Getachew and D. Jabesa, "Afaan Oromo Text Content-Based Fake News Detection using Multinomial Naive Bayes," International Journal of Innovations in Management, Science and Engineering (IJIMSE), vol. 01, no. 01, pp. 26-37, 01 March 2020.

- [7] I. y. R. Pratiwi, "study of hoax detection using neive bayes classfier in indonesian language," in International Conference on Information & Communication Technology and System (ICTS), 2017.
- [8] D. Pisarevskaya, "Deception Detection in News Reports in the Russian Language," 2019.
- [9] N. J. Conroy, "Automatic deception detection: Methods for finding fake news,," in in Proceedings of the 78th ASIS&T Annual Meeting: Information Science with Impact: Research in and for the Community, USA, 2015.
- [10] J. D'Souza, "An Introduction to Bag-of-Words in NLP," 03 04 2018. [Online]. Available: https://medium.com/greyatom/an-introduction-tobag-of-words-in-nlp-ac967d43b428.
- [11] G. Bonaccorso, "Artificial Intelligence Machine Learning – Data Science," 10 06 2017. [Online]. Available: https://www.bonaccorso.eu/2017/10/06/mlalgorithms-addendum-passive-aggressivealgorithms/.
- [12] Gilda, "Evaluating Machine Learning Algorithms for Fake News Detection," in IEEE 15th Student Conference on Research and Development (SCOReD), 2017.