

Fine-Grained Fake News Detection System

Harsh Goyal
harsh20562@iiitd.ac.in

Ayush Sharma
ayush20042@iiitd.ac.in

Anupam Narayan
anupam20030@iiitd.ac.in

Shivam Jindal
shivam20125@iiitd.ac.in

Anas Ahmad
anas20023@iiitd.ac.in

April 23, 2023

1 Motivation

False or misleading information presented as news is known as fake news. Most of the time, fake news aims to hurt someone or something's reputation or make money from advertising. Of late, fake news has been in the spotlight of mainstream journalism and the general public because it can affect a country's political scenario. In this project, we attempt to detect the authenticity of specifically political news. Social media is the primary channel for disseminating such content, though it occasionally becomes mainstream media. Because fake news can significantly impact an election's political outcome, it is becoming increasingly important to identify and classify it as such. The loose definition of "fake news" presents the primary obstacle to resolving the problem. For instance, fake news can be broken down into subcategories: a statement known to be completely false, a speech that presents statistics as facts that have not been thoroughly investigated, or satirical writing. Our main aim of the project is to detect and classify fake news.

2 Problem Statement

Fake news is false or misleading information presented as news. Multiple strategies for fighting fake news are currently being actively researched for various types of fake news. Our task is to determine whether a statement/speech spoken by a politician in a political rally and when published on social media/news articles as text is fake or not using NLP and associated language models. We aim to design novel and hybrid models using

pre-trained large language models to integrate the original news statement/speech with the metadata. Fake News is a text-classification task. Since the news statements are very short in length and the text from the speech is noisy and contains grammatical errors, the task becomes more complex and exciting.

3 Literature review

3.1 Liar, Liar Pants on Fire

The author proposed a hybrid Convolution Neural Networks Framework for integrating text and metadata. Hybrid CNN consists of two parallel Convolution Layers in which the input to the first Convolution Layer is the word embeddings for the given statement, followed by the Max-Pooling. In contrast, the input to the second convolution layer is the metadata like speaker name, subjects of the speech./statement, speaker's party, etc., followed by a Bidirectional-LSTM layer. The output of both layers is then concatenated and fed into a fully-connected layer with softmax on the output layer.

3.2 A Retrospective Analysis of the Fake News Challenge Stance Detection Task

The author uses a fake news challenge dataset which is a 4-class classification task. The author used two stacked LSTMs with 50-dimensional GloVe word embeddings as input and a three-layered neural network with 600 neurons each. The output from a three-layered neural network is fed into the output layer consisting of 4 neurons to estimate the class-wise probabilities to estimate to which

class it belongs.

3.3 Exploring Text-transformers in AAI 2021 Shared Task: COVID-19 Fake News Detection in English

The author proposed two solutions to the problem. The first solution uses RNNs (also called Bidirectional-LSTMs) as the LSTMs are based on previous text information. In contrast, the RNNs are based on both previous and later text information, which helps get a better sentence context. The second solution uses 3 different techniques - a Five-Fold Single-model ensemble, a Five-Fold Five-model ensemble, and a Pseudo Label algorithm. In the Five-Fold Single-model ensemble, the author used the same transformer-based pre-trained models (like BERT, Roberta, etc.) on all the five-folds of the dataset. The author used different pre-trained models for fine-tuning each fold in the Five-Fold Five-model ensemble. In the Pseudo Label algorithm, the author proposed that if the test data sample is classified to some class with a probability greater than 0.95, then the author proposed to use that test data sample as the training data sample for future test samples.

3.4 On the Benefit of Combining Neural, Statistical and External Features for Fake News Identification

The author proposed how using three different word embeddings in parallel and concatenating them can benefit fake news detection. The author combined neural embeddings, statistical features, and external features to a similar type of model (with a little difference in activation functions) in parallel and concatenated them to get the combined features. The author defined neural embeddings as the skip-thought vectors which encode the given sentences to vector embedding of length 4800. The statistical features are defined as the vectors obtained from the text using Bag-of-words(BOW), TF-IDF, and n-grams techniques. The external features include heuristics such as the similarity between headline-body pairs, the number of similar words in the headline and the body, etc. All these three types of features are combined using pre-trained models and then fed into dense layers to predict

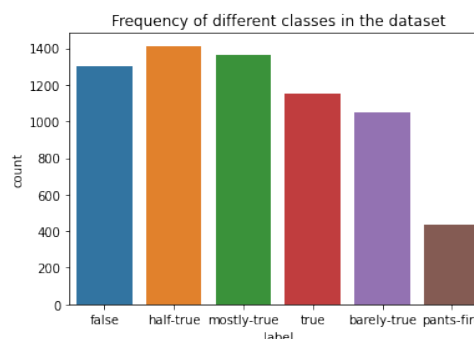
the correct label for the given input sentence.

4 Dataset Details

Dataset Statistics	
Training set size	10,269
Validation set size	1,284
Testing set size	1,283
Avg. statement length(tokens)	17.9

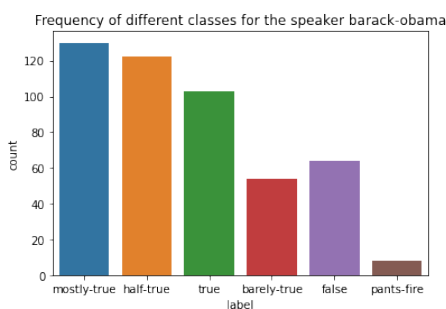
The 6-classes of the dataset are defined as:- (These definitions were taken from PolitiFact's "truth-o-meter" methodology page).

1. **true** – The statement is accurate, and there's nothing significant missing.
2. **mostly-true** – The statement is accurate but needs clarification or additional information.
3. **half-true** – The statement is partially accurate but leaves out important details or takes things out of context.
4. **barely-true** – The statement contains an element of truth but ignores critical facts that would give a different impression.
5. **false** – The statement is not accurate.
6. **pants-fire** – The statement is inaccurate and makes a ridiculous claim. a.k.a. "Liar, Liar, Pants on Fire!"



4.1 Updated Dataset

As the training dataset is minimal (after removing NaN rows from the dataset, the amount of data decreased further), therefore to generalize our models better, we merged the validation data in the training data and tested our models on the testing data.



5 Models and Methodology

5.1 Baseline Models

We computed the TFIDF matrix for the text statements in the training data. For training purposes, we gave TFIDF matrix and the corresponding labels as input to the classification model and then tested the model for the testing data. We tried this technique on various ML models like Logistic Regression, Naive Bayes, Random Forest, etc. We also tried hyperparameter tuning for all these ML models.

5.2 Mid-Review Evaluation Models

We computed POS taggings for the given text data. We used TFIDF on the POS taggings to convert textual data to numerical data and fed this data to various machine learning models like Logistic Regression, Multi-Layer Perceptron, etc. We computed non-contextual word embeddings using GloVe embeddings which converts textual data to 100 dimensional vector. We also computed contextual embeddings using sentence transformer model all-MiniLM-L6-v2 which converts textual data to 384 dimen-

sional vector for each token based on the context of the sentence. We gave non-contextual and contextual word embeddings as input to various machine learning models.

5.3 Recurrent Neural Networks

We first tokenize the statements using the BERT tokenizer and train an RNN-based BERT model to predict the labels. The model consists of a BERT layer followed by an RNN layer. The RNN part of the code is responsible for processing the output of the BERT layer and producing a final output for the model. Specifically, the RNN layer used in this code is a Bidirectional LSTM (Long Short-Term Memory) layer. The model is trained for three epochs using a batch size of 32, and the test accuracy is reported at the end.

5.4 CNNs / LSTM / BiLSTM

For all these models, we created POS encodings for the preprocessed data. We passed the textual data and POS encodings to Embedding() layer to get the embeddings and passed these embeddings as input to the these models. We used categorical cross-entropy as a loss function. We ran these models using two approaches - only giving embeddings of statement as input and giving embeddings of text and POS encodings as input to these models.

5.4.1 LSTM

- LSTM size = 100
- hidden-size = 100
- num-epochs = 30
- batch-size = 40
- embedding-dim = 100
- dropout = 0.2
- loss function = categorical_crossentropy
- num-steps = 15 : words in a stmt
- optimizer = SGD(lr=0.025, clipvalue=0.3, nesterv=True)

5.4.2 CNNs

- kernel sizes for each of conv1d = 3
- num-epochs = 30
- filter-size for each of conv1d = 128
- batch-size = 40

- dropout = 0.2

- loss function = categorical crossentropy

- optimizer = SGD(lr=0.025, clipvalue=0.3, nesterv=True)

5.5 Pre-trained models (BERT, RoBERTa, XLNet)

For bert we used 'bert-base-uncased' model, for RoBERTa we used 'roberta-base' and for XLNet we used xlnet-base-cased. All these models are pretrained models. We first need to preprocess textual data into numbers which can be done by inbuilt tokenizer of each of these models. We then need to train the model on the input data and labels and then we evaluate the model on validation set after each epoch and finally on the testing set after complete training. We used 3 epochs, batch size of 32 and learning rate as $2e-5$.

6 Evaluation Metric

For the baseline results, earlier, we used weighted F1-score as our evaluation metric. But, after doing some research and based on our understanding of the concepts, we thought that the weighted F1-score was not a good metric for the given dataset. According to our understanding, the weighted F1-score is generally used when the number of samples corresponding to each class in the dataset is similar, as the weighted F1-score gives more weightage to more regular classes as compared to less frequent classes while computing the weighted F1-score. So, we decided to change our evaluation metric from weighted F1-score to **macro F1-score** as the macro F1-score gives equal weightage to all the classes, unlike the weighted F1-score. For the final submission, we also used **accuracy** as one of the evaluation metrics.

7 Experimental Results

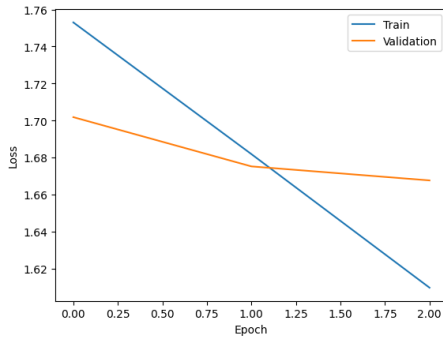
Models	F1-Score	Accuracy
CV + LR	0.22	0.25
TF-IDF + SVM	0.22	0.26
CV + DT	0.22	0.25
Other Baseline Models	0.18-0.21	0.21-0.25
CNN(Statements)	0.21	0.23
CNN(Statements + POS)	0.17	0.23
RNN	0.26	0.26
LSTM(Statements)	0.20	0.22
LSTM(Statements + POS)	0.21	0.22
BiLSTM(Statements)	0.21	0.25
BiLSTM(Statements + POS)	0.21	0.22
Bert	0.278	0.283
Roberta	0.286	0.292
XLNet	0.279	0.289

The results presented show the performance of various models on the given dataset. The evaluation metric used is the macro F1 score and accuracy.

Looking at the results, we beat State of the Art accuracy (**0.274**) using pre-trained RoBERTa (**0.292**) (We merged validation data into the training data as training data was pretty small compared to number of features.)

The highest-performing model is the pre-trained RoBERTa model, which also outperformed the pre-trained BERT and XLNet models. This might be due to the fact that RoBERTa uses dynamic masking in each epoch while BERT uses static masking in each epoch and hence, RoBERTa is generalized better. This indicates that this model achieved the best balance between precision and recall for both the positive and negative classes.

The other models like BERT and XLNet outperformed other deep learning models like CNNs, LSTMs, BiLSTMs, etc. Apart from the pre-trained models, Recurrent Neural Networks performed best as compared to other Deep Learning and Machine Learning models.



8 Future Work

For future work, we will try to implement Hybrid Siamese Network with the following approaches:

- **Hybrid Input:** In this approach, we give text statements as input to one BERT and metadata as input to another BERT, and the output is concatenated and fed into a dense layer.
- **Hybrid Pre-trained models:** In this approach, we plan to use different pre-trained models as the base models for the Siamese Network.

9 References

- [1] William Yang Wang. 2017. “Liar, Liar Pants on Fire”: A New Benchmark Dataset for Fake News Detection. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 422–426, Vancouver, Canada. Association for Computational Linguistics.
- [2] Andreas Hanselowski, Avinesh PVS, Benjamin Schiller, Felix Caspelherr, Debanjan Chaudhuri, Christian M. Meyer, and Iryna Gurevych. 2018. A Retrospective Analysis of the Fake News Challenge Stance-Detection Task. In Proceedings of the 27th International Conference on Computational Linguistics, pages 1859–1874, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- [3] Li, X., Xia, Y., Long, X., Li, Z., and Li, S. (2021). Exploring Text-transformers in AACL 2021 Shared Task:

COVID-19 Fake News Detection in English. ArXiv. <https://doi.org/10.48550/arXiv.2101.02359>

[4] Bhatt, G., Sharma, A., Sharma, S., Nagpal, A., Raman, B., and Mittal, A. (2017). On the Benefit of Combining Neural, Statistical and External Features for Fake News Identification. ArXiv. <https://doi.org/10.48550/arXiv.1712.03935>

10 Contributions

Motivation - Anas Ahmad, Shivam Jindal

Problem Statement - Anupam Narayan, Ayush Sharma

Literature Review - Harsh Goyal

Models and Evaluations - Harsh Goyal, Ayush Sharma, Shivam Jindal, Anupam Narayan