



# Comparative Analysis of Machine Learning Algorithms for Fake News Detection

Nikhil Madaan  
Siddharth Dharm  
Prof. Aruna Malapati

# Introduction

- In the IoT era, no less than twenty billion devices are connected to the Internet, and the dissemination of information also has become inexpensive, hassle free.
- Moreover, the internet has made it easy to post content without any restrictions.
- The propagation of news (both, legitimate and fake) is inevitable.
- Several research groups have delved into this problem.
- Very few of them have made use of Linguistic properties of the English (natural) language.

# Dataset - FakeNewsCorpus

- This dataset contains 9,408,908 labelled news articles. These articles have been scraped from a curated list of 1001 domains from <http://www.opensources.co/>.
- This corpus is aimed toward training algorithms for detecting fake news. Includes
  - news articles related to a number of 'tags', like fake, reliable, satire, bias, etc.
- 60,000 articles were selected randomly from the corpus.
  - Articles had their 'tag' attribute either equal to 'fake' or 'reliable'.
- During random sampling from the corpus, it was ensured that articles belonging to both the categories, were represented equally.

# Features

- Word Embeddings
  - The text is being represented using word embeddings (Mikolov et al. 2013a). We have used a pre-trained Google word2vec model (Mikolov et al. 2013b) to get the vectorized representations (of dimensions) for the words.
  - represent the text, the mean of Word Embeddings are taken into account.
- Syllable Count - Returns a count of the number of syllables in a particular text.
- Sentence Count - Returns a count of the number of sentences in a particular text.
- Flesch-Kincaid Readability Tests - Represents the ease with which a particular text can be read. It provides the following metrics :-
  - Flesch-Kincaid Grade Level
  - Flesch Reading Ease

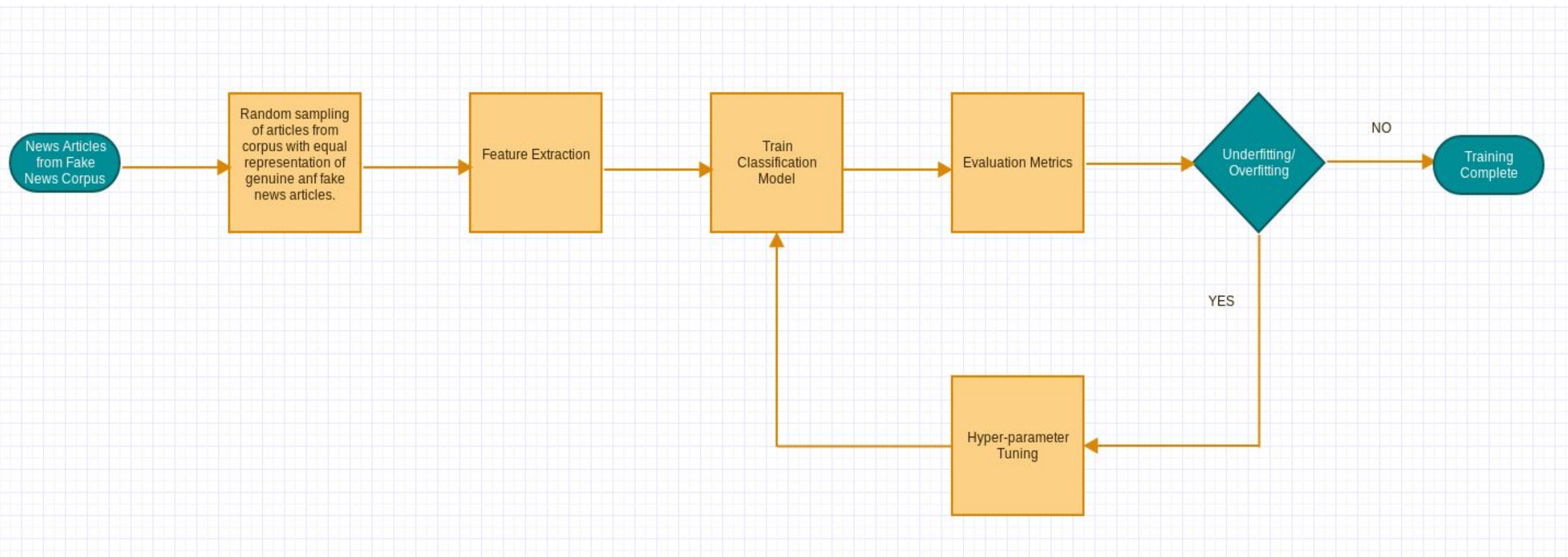
# Features Contd.

- Gunning-Fog Index
  - Produces an estimate of the number of years of formal education required by a person to understand a given text on the first reading.
  - Different from Flesch - Kincaid Grade Level.
- Automated Readability Index
  - Gauges the understandability of a text. Given by (Eltorai et al. 2015)
- SMOG Grade
  - Stands for 'Simple Measure of Gobbledygook'.
  - Returns a measure of readability similar to the Gunning Fog Index that estimates the years of education needed to understand a piece of writing.
- Linsear Write Score
  - This score is designed to calculate the readability of technical writing.

# Features Contd.

- Dale Chall Readability Score
  - Provides a numeric gauge of the difficulty of comprehension that readers face while reading a given text.
  - Formula given by (Dubay 2004)
- Coleman-Liau Index
  - Used to gauge the understandability of a text.
  - Given by (Karmakar and Zhu 2010).
- Part-of-Speech (POS) Tagging
  - We restrict the POS features to Nouns, Adjectives, Verbs and Adverbs since they are most informative and generic unlike other parts of speech, such as Conjunctions and Prepositions.

# Workflow



# Models Compared

- Logistic Regression
- Random Forest
- Support Vector Machines
- Artificial Neural Network
- LSTM (Long Short-Term Memory)
  - Allows us to implicitly capture temporal structure in the sentence.
- Bi-directional LSTM (Long Short-Term Memory)
  - Provides capacity to process of sentence in both forward and backward direction.



# Models Compared

- GRU (Gated Recurrent Unit)
  - Another instance from family of RNNs with comparable performance to LSTMs but with lesser parameters.
- Bi-directional GRU (Gated Recurrent Unit)
  - To provide bi-directional information processing capacity
- Our work demonstrates a comparative analysis of the results obtained by setting the dimensions of word embeddings as 50, 100 and 200.

# Results

**Table 1: Evaluation Metrics of Classification Algorithms**

Algorithm	Accuracy	Sensitivity	Specificity	Precision
Logistic Regression	0.9444	0.9397	0.9492	0.9493
Random Forest	0.9481	0.9302	0.9659	0.9645
SVM	0.8947	0.8836	0.9045	0.9026
ANN	0.9574	0.8636	0.9454	0.9416
LSTM(50-D)	0.9721	0.9866	0.9574	0.9587
Bi-LSTM(50-D)	0.9796	0.9895	0.9774	0.9777
LSTM(100-D)	0.9723	0.9682	0.97627	0.97620
Bi-LSTM(100-D)	0.9765	0.9875	0.9654	0.9663
LSTM(200-D)	0.9669	0.9945	0.9388	0.9429
Bi-LSTM(200-D)	0.9828	0.9945	0.9388	0.9429
GRU(50-D)	0.9312	0.9895	0.9249	0.9347
Bi-GRU(50-D)	0.9290	0.9857	0.9380	0.9355
GRU(100-D)	0.9225	0.9790	0.9372	0.9343
Bi-GRU(100-D)	0.9232	0.9744	0.9315	0.9389
GRU(200-D)	0.9222	0.9667	0.9318	0.9451
Bi-GRU(200-D)	0.9258	0.9704	0.9446	0.9457

# Conclusion

- In this paper, we proposed a deep learning solution to the problem and have presented a comparative analysis of how standard machine learning algorithms fare against recurrent neural networks.
- We observed that Bi-LSTM performed better than other algorithms with reasonable performance gains as can be deduced from the values of evaluation metrics obtained.