

# Fake News Detection

Using CNN-LSTM Hybrid Model Approach

**Anirudh Bagalkotker - 2021A7PS2682H**

**Kartik Pandey - 2021A7PS2574H**

**Research Paper:**

**<https://ieeexplore.ieee.org/abstract/document/9178321>**



# Abstract

The rise of fake news presents serious challenges to both political stability and social trust, as misinformation spreads rapidly across digital platforms. In an era dominated by social media, the ability to quickly and accurately assess the legitimacy of news has become increasingly difficult. To address this challenge, we propose an automated fake news detection model using a **hybrid neural network architecture** that integrates **Convolutional Neural Networks (CNN)** and **Long Short-Term Memory (LSTM) networks**. This model leverages a Hybrid approach, to **enhance the classification performance** by reducing the complexity of feature vectors before passing them through the neural network.

The model is will be trained on the Fake News Dataset from Kaggle, which categorizes news articles into four stances: agree, disagree, discuss, and unrelated. This work highlights the potential of combining deep learning models with dimensionality reduction to create more reliable and efficient tools for detecting fake news.

# Dataset & Features

This dataset contains both fake and real news articles, aimed at helping distinguish between genuine and false information presented as news. It's designed to assist in the study and detection of fake news, which is often created to damage reputations or generate advertising revenue. The size of the data is 6256.

## Features:

- 1. Title:** This column contains the headline or title of the news article. Headlines are often designed to capture attention and may play a significant role in distinguishing real news from fake news, as misleading or exaggerated titles are common in fabricated stories.
- 2. Text:** This feature contains the full body of the news article. The content of the article provides crucial information, as fake news often includes deceptive or incomplete facts, which can be detected through text analysis.
- 3. Label:** The label assigns a categorical value to each news article—either "REAL" or "FAKE." This column serves as the target variable for the classification model and is used to train the machine learning algorithm to differentiate between genuine and fabricated news.

Source: <https://www.kaggle.com/datasets/rajatkumar30/fake-news>

# Dataset & Features

This dataset would be split into 3 Segments using stratified random splitting with a train-validation-test split.

- 1. Training Set (60%):** Used for training the model.
- 2. Validation Set (20%):** Used to fine-tune hyperparameters and evaluate the model's performance during training, allowing you to avoid overfitting.
- 3. Test Set (20%):** Held out until the very end, used to assess the model's final performance on unseen data.

The data is split in a way that preserves the proportion of each class label ("REAL" or "FAKE") across the train, validation, and test sets. This ensures that each set maintains a similar class distribution to the original dataset, which is especially important if your data is imbalanced.

This type of split is common in machine learning tasks to ensure robust model development and evaluation, preventing overfitting and providing reliable performance metrics.

# Methodology

The proposed fake news detection system is built using a hybrid deep learning model, combining Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks. The following outlines the key methodologies, from data preparation to model evaluation.

## Data Preprocessing:

- 1. Text Cleaning:** Both the title and the body text of the articles are preprocessed by removing punctuation, special characters, stopwords, and converting all text to lowercase to ensure uniformity.
- 2. Tokenization:** The text is tokenized into individual words, splitting the sentences into smaller units for analysis.
- 3. Word Embeddings:** Word embeddings like Word2Vec or GloVe are used to convert tokens into dense vectors that represent the semantic meaning of words, enabling the model to understand the context.

A CNN is used to capture local patterns in the text, while LSTM models sequential dependencies. The outputs are combined and passed through fully connected layers to classify news as REAL or FAKE.

# Methodology

The hybrid model architecture consists of:

- 1. CNN:** Convolutional layers extract local patterns in the text, capturing essential n-gram features that may signal fake or real news.
- 2. LSTM:** This recurrent network is applied to capture the temporal dependencies and contextual information in the news text, enabling the model to understand word sequences.

The outputs from CNN and LSTM are merged and passed through fully connected layers, with a final softmax layer to classify the news as REAL or FAKE.

The training set is used to fit the model, while the validation set assists in fine-tuning and hyperparameter optimization.

The model's performance is evaluated using standard metrics such as accuracy, precision, recall, and F1-score. These metrics provide a comprehensive view of the model's effectiveness in distinguishing between real and fake news.