

Predicting news truthfulness through graph-based retweet patterns.

Baggio Davide 2122547
Martinez Zoren 2123873
Brocheton Damien 2133034

Motivation

The rise of misinformation on social media has significant implications for public opinion, health, and safety, making it crucial to distinguish real news from fake. Twitter, as a major news source, often spreads information rapidly, sometimes without verification. By analyzing the graph structure of news propagation on Twitter, we can identify patterns in how real and fake news spread. This project aims to develop insights and tools to enhance the credibility of online information, contributing to a more informed and resilient public. Recent studies[1] have shown that machine learning models can effectively detect real or fake news by analyzing user-specific data, such as the profiles of those sharing the information. One of the biggest catches is the complexity of the generated models, mostly being Convolutional Neural Networks applied to Graphs and the low accuracy given new data. In this project, however, we aim to explore whether it's possible to classify news as real or fake based solely on the "pure" retweet graph structure, independent of user metadata. By employing the algorithms outlined in the following sections, we will extract essential features from the retweet graph that can serve as inputs for a machine learning model, enabling an analysis based purely on the patterns of information spread.

Dataset

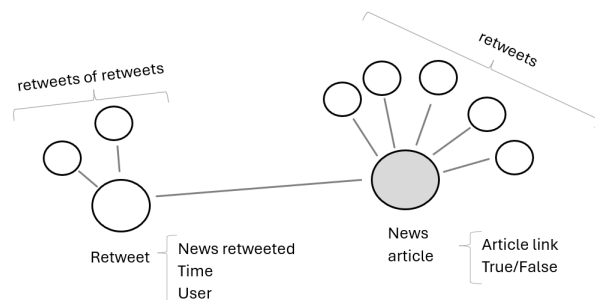
The dataset is part of a bigger pool provided by the Twitter API. This part is shared on github under the Apache Licence, Version 2.0[2]. The dataset is well documented in this paper[3]. It is basically a Graph with many connected components, each representing a news tree, composed of the main tweet of the news and all the retweets associated with it.

Method

Problem

Our objective is to identify the characteristics of tweets that reference a fake news, and determine if a new tweet references it by looking at its characteristics. To analyze the dataset, we need to construct the graph from the dataset. The structure is represented as follows:

- **Graph Type:** Tree-structured graphs
- **Root Node:** News item (labeled true/false)
- **Other Nodes:** Twitter users who retweeted the news
- **Edges:**
 - News item \leftrightarrow User: Direct retweet
 - User \leftrightarrow User: Retweet through an intermediary
- **Additional Information:** Retweet timestamps



The goal is to extract features from the network and use that data to train a machine learning model to predict the truthfulness of news. The algorithms that we are going to use are important for:

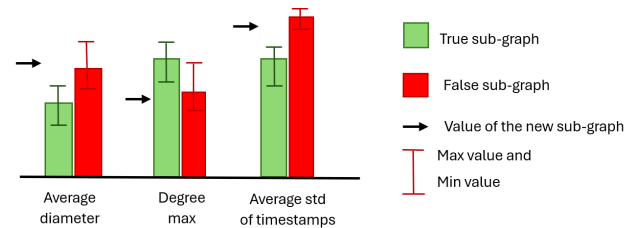
- **Average depth of trees:** Is the depth of the trees greater for fake news?
- **Average retweet breadth:** Do fake news tend to spread more quickly having a greater breadth at the first level?
- **Average time between retweets:** Do fake news spread faster?
- **Peak diffusion time:** Do fake news have an explosive peak?
- **Users reliability score:** various features based on the ranking of users who shared the news
- **Centrality and Pagerank:** What nodes are more common to find in paths between other 2 nodes?

Intended Experiments

User reliability ranking

One of the goal is to visualize wether new data involving a news fits better into the "Real news" or "Fake news" category. This can be achieve displaying a bar graph representing the features studied previously.

The final step is to create a ranking of users based on their reliability and to generate features for the data using this ranking. To do this, the following factors will be considered:



- Number of retweets of false news
- Percentage of retweets of true vs. false news
- Betweenness centrality of the user in the graph
- Pagerank of the user in the graph

For large datasets computing the ranking for ALL users might be computationally expensive. Therefore, in this case, users with good centrality measures will have priority. The features of the news obtained from this ranking are as follows:

- Average Ranking Positions of users who retweeted the news
- Percentage of reliable users (e.g. percentage of users in the top 20)
- Minimum and maximum position
- Weighted score based on reliability

Machine learning models

The models we want to try are the following:

- **Support Vector Machine**
- **Feed Forward Neural Network**
- **Random Forest**

We will make comparisons to better understand which model fits this problem best, based on which one achieves better accuracy and also the time it will take to train it.

Libraries: Networkx[4] (for Graph analysis), Scikit-Learn (for SVM and RF models), Tensorflow (for FFNN model)

Evaluation metrics of the model: Accuracy, Precision, Recall and F1-Score

Machine for experiments:

- AMD Ryzen 5 3500U (8-cores), 8GB DDR4, Windows 11 or Ubuntu 20.04
- AMD Ryzen 5 4500 (8-cores), Radeon RX 6600, 16GB DDR4, Windows 11 or Ubuntu 20.04

Development

Extracting informations from the dataset

The first part of the project involves extracting the data from the dataset. As shown previously the dataset is composed of many connected components, each representing a news (related to gossip or politics). It has been decided, during the development, to consider only few features that could be relevant to the objective of the project. Those features are:

Diameter: provides insight into the maximum extent of information spread. The bigger the diameter, the deeper (given that the structure of the graph is a tree) the news is spreaded among the users.

Maximum degree: identifies the node with the highest number of connections. In the case of fake news, such nodes could be targeted for spreading misinformation widely.

Degree centrality: the number of direct connections a node has. Nodes with high degree centrality are crucial in the immediate dissemination of news.

Closeness centrality: how quickly information can spread from a given node to all other nodes in the network. High closeness centrality suggests that a node is well-positioned to efficiently spread news.

PageRank: rank of nodes in a graph based on their importance. High PageRank nodes are influential and could be key disseminators of real or fake news.

Average Standard Deviation of Retweet Timestamps: consistency or burstiness of retweet activity over time. A small average standard deviation means that the news could be spreaded in a small amount of time since it has been tweeted for the first time.

Using networkx as the library for analyzing the graphs, extracting these information was actually easy. With all the features we then exported them into ".csv" files and ultimately we calculated all the average values separately for news labeled real and fake that are going to be used later in a system of prediction based on score.

```
1  import networkx as nx
2
3  #calculating std of the timestamps
4  std = np.std(np.array(timestamps))
5  #calculating the diameter
6  d = nx.diameter(s.graph)
7  #calculating the max degree
8  _, neighbors = max(s.graph.degree, key=lambda x: x[1])
9  #calculating degree centrality
10 dc = np.mean(list(nx.degree_centrality(s.graph).values()))
11 #calculating closeness centrality
12 cc = np.mean(list(nx.closeness_centrality(s.graph).values()))
13 #calculating pagerank
14 pr = np.mean(list(nx.pagerank(s.graph).values()))
15
16 #function for exporting data into .csv
17 def dump_data(data: list[list]):
18     file_name = f"./{file_type}.csv"
19     with open(file_name, "w") as file:
20         for graph in data:
21             s = ""
22             for value in graph:
23                 s += str(value) + ", "
24             s = s[:-2] + "\n"
25             #print(s)
26             file.write(s)
```

References

- [1] Yi Han, Shanika Karunasekera, Christopher Leckie
"Graph Neural Networks with Continual Learning for Fake News Detection from Social Media",
<https://arxiv.org/pdf/2007.03316>, 2020.
- [2] Dataset: <https://github.com/safe-graph/GNN-FakeNews/tree/main>,
https://drive.google.com/drive/folders/10slTX91kLEYIi2WBnwuFtXsVz5SS_XeR?usp=sharing
- [3] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee and Huan Liu
"FakeNewsNet: A Data Repository with News Content, Social Context and Spatiotemporal Information for Studying Fake News on Social Media",
<https://arxiv.org/pdf/1809.01286>, 2019
- [4] NetworkX library: <https://networkx.org/documentation/stable/>

Contribution

Contributors:

- Baggio Davide ($\frac{1}{3}$ of the work): Finding a well documented dataset, reading the related papers and understanding it. Writing the first part of the proposal and preprocessing the dataset ready for analysis into a python script.
- Martinez Zoren ($\frac{1}{3}$ of the work): Finding a well documented dataset, reading the related papers and understanding it. Writing the third part of the proposal and planning on the machine learning models to use in order to achieve the goal of the project.
- Brocheton Damien ($\frac{1}{3}$ of the work): Finding a well documented dataset, reading the related papers and understanding it. Writing the second part of the proposal and starting to write the post-processing of the data into a python script that compare new data with the studied one using a probabilistic algorithm.