

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

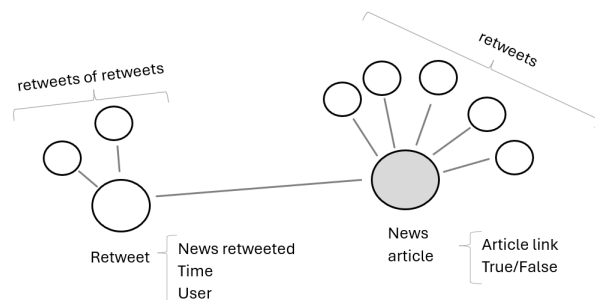
The dataset is part of a bigger pool provided by the Twitter API. This part is shared on github under the Apache Licence, Version 2.0[2]. The dataset is well documented in this paper[3]. It is basically a Graph with many connected components, each representing a news tree, composed of the main tweet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

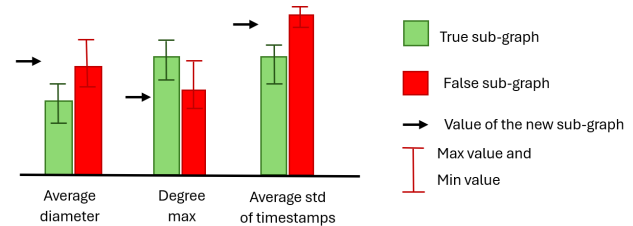
- **Average depth of trees:** Is the depth of the trees greater for fake news?
- **Average retweet breadth:** Do fake news tend to spread more quickly having a greater breadth at the first level?
- **Average time between retweets:** Do fake news spread faster?
- **Peak diffusion time:** Do fake news 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[5] (for SVM and RF models), Tensorflow[6] (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[4] 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()))
```

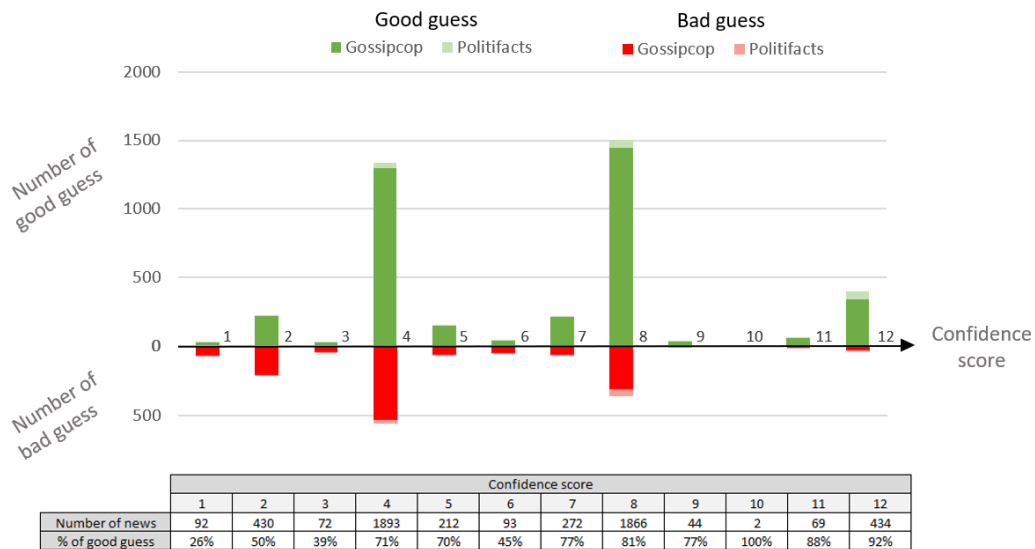
Prediction system

Once the features are extracted we have created a system based on score that is capable of predicting in a probabilistic way if a news is real or fake. This system gets some graphs data as input and determines how much it deviates from the averages calculated with previous algorithms giving a final score from 0 to 12 (12 being the highest confidence).

```

1  def confidence_scoring(value, avg_r, avg_f):
2      # Compute the difference
3      diff_r = abs(avg_r - value)    #avg of real sub-graph - new value
4      diff_f = abs(avg_f - value)    #avg of fake sub-graph - new value
5      inside = (value < avg_r and value > avg_f) or (value < avg_f and value > avg_r) #The
        value is between the 2 avg
6      #Closest one
7      closer = 0                    # 0 = Equal
8      if diff_r < diff_f:
9          closer = 1                # 1 = Real
10     elif diff_r > diff_f:
11         closer = 2                # 2 = Fake
12     #Confidence score
13     score = 0                      #No guess
14     if closer == 1:                #Real
15         if 2*diff_r < diff_f or not inside:
16             score = 2              #Strong guess
17         else:
18             score = 1              #Weak guess
19     elif closer == 2:              #Fake
20         if 2*diff_f < diff_r or not inside:
21             score = -2             #Strong guess
22         else:
23             score = -1             #Weak guess
24     return score

```



The image above is just an example on how it can display the prediction and its accuracy. It represent a confidence score on how sure we are about a news label. From what we can see, there is a correlation between the data extracted from the graphs and their actual label. This is what we expected from the beginning so this will be important to develop a machine learning model which is able to predict the truthfulness of a news with a sufficiently high accuracy.

Machine learning models and User reliability ranking

In the final part of the project some ML models have been developed and tested onto the features retrieved before. By using the dataset we were able to map the nodes ID to the user ID and determine the number of real and fake news retweeted by the users. We then calculated the user score with this equation:

$$\begin{aligned} \text{UserScore} &= \text{US} \\ \text{Normalized number of false retweets} &= \text{FR} \\ \text{Normalized true to false ratio} &= \text{TFR} \\ \text{Normalized pagerank of node} &= \text{NP} \\ \text{Normalized degree centrality of node} &= \text{NP} \\ \text{Normalized closeness centrality of node} &= \text{NP} \end{aligned}$$

$$US = 0.4 \times FR + 0.3 \times TFR + 0.2 \times NP + 0.1 \times DC + 0.1 \times CC \quad (1)$$

```
1  # Calculate the score for each user
2  # Calculate FR (False Retweets)
3  FR = user["FalseNews"] / user["TotalNews"]
4      if user["TotalNews"] > 0 else 0
5  # Calculate TFR (True to False Ratio)
6  TFR = user["TrueNews"] / user["FalseNews"]
7      if user["FalseNews"] > 0 else user["TrueNews"]
8  # Calculate NP (PageRank)
9  NP = (user["PageRankCentrality"] - min_pr) / (max_pr - min_pr)
10     if max_pr != min_pr else 0
11 # Calculate DC (Degree Centrality)
12 DC = (user["DegreeCentrality"] - min_dc) / (max_dc - min_dc)
13     if max_dc != min_dc else 0
14 # Calculate CC (Closeness Centrality)
15 CC = (user["ClosenessCentrality"] - min_cc) / (max_cc - min_cc)
16     if max_cc != min_cc else 0
17 # Penalize users with low TotalNews
18 penalty = alpha / user["TotalNews"] if user["TotalNews"] > 0 else 0
19 # Calculate UserScore
20 user["UserScore"] = 0.4 * FR + 0.3 * TFR + 0.1 * NP + 0.1 * DC + 0.1 * CC
```

This user ranking will be useful later to train some models in order to increase the general accuracy. As stated before we are going to implement three type of ML models: Support Vector Machine, Feed Forward Neural Network and Random Forest of which all have been tuned with a validation set to have the best parameters. The final question is: how well does these model perform on the features extracted before?

Results

```
Test accuracy: 0.746031746031746
Classification Report:
      precision    recall  f1-score   support

     0.0         0.81     0.71     0.76         35
     1.0         0.69     0.79     0.73         28

 accuracy         0.75
 macro avg         0.75
 weighted avg      0.75
```

(a) Random Forest Results

```
GENERALIZATION SCORE BEST MODEL
35/35 0s 4ms/step - accuracy: 0.8490 - loss: 0.3586
Test accuracy: 0.8536139130592346
35/35 0s 3ms/step
Classification Report:
      precision    recall  f1-score   support

     0.0         0.85     0.86     0.86         561
     1.0         0.85     0.85     0.85         532

 accuracy         0.85
 macro avg         0.85
 weighted avg      0.85
```

(b) Feed Forward Neural Network Results

```
Training data size: 4370
Testing data size: 1094
Training size: 4370
Test size 1094

Training Poly SVM with C=10 and degree=3
Training score: 0.8368421052631579

Generalization (Test) score: 0.8244972577696527

Classification Report:
      precision    recall  f1-score   support

     0.0         0.84     0.80     0.82         547
     1.0         0.81     0.85     0.83         547

 accuracy         0.82
 macro avg         0.83
 weighted avg      0.83
```

(c) Support Vector Machine Results

The Support Vector Machine results were promising even though the model is very simple, requires few resources and the training is not time consuming. We then tried to concatenate the User ranking data to the features of the news graph and feed them into another SVM model with the same parameters as the one used in figure (a). The results were impressive:

```
Training main data size: 4370
Training additional data size: 4370
Testing main data size: 1094
Testing additional data size: 1094
Training size: 4370
Test size 1094

Training Poly SVM with C=10 and degree=3
Training score: 0.9956521739130435

Generalization (Test) score: 0.9204753199268738

Classification Report:
      precision    recall  f1-score   support

     0.0         0.96     0.88     0.92         547
     1.0         0.89     0.96     0.92         547

 accuracy         0.92
 macro avg         0.92
 weighted avg      0.92
```

(d) Support Vector Machine with ranking Results

Issues encountered

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/>
- [5] Scikit-learn library: <https://scikit-learn.org/stable/api/index.html>
- [6] Tensorflow library: https://www.tensorflow.org/api_docs/python/tf

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.