

Named Entity Recognition and Its Effects on News Article Engagement

1. Introduction

The proliferation of fake news and its impact on public perception has become a critical area of study. This report investigates the relationship between named entities (such as organizations, people, and geopolitical entities) in news article titles and their popularity (measured through engagement on Twitter). The analysis also explores the correlation between named entities and sentiment with the real/fake nature of news, providing actionable insights for detecting fake news patterns.

2. Methodology

2.1 Datasets

The analysis utilized four datasets:

1. **GossipCop Fake**
2. **GossipCop Real**
3. **Politifact Fake**
4. **Politifact Real**

Each dataset contained the following columns:

- **id**: Unique identifier for each article.
- **news_url**: The URL of the article.
- **title**: The headline of the article.
- **tweet_ids**: A comma-separated list of tweet IDs referencing the article.

To combine these datasets, additional columns were introduced:

- **label**: Indicates whether the article is "real" or "fake".
- **source**: Identifies the dataset source ("GossipCop" or "Politifact").

2.2 Data Preprocessing

- Headlines were cleaned (converted to lowercase and stripped of extra spaces) to ensure consistent analysis.
- Named Entity Recognition (NER) was performed using the SpaCy model `en_core_web_sm`. The entities tracked included:
 - **ORG**: Organizations.
 - **PERSON**: People.
 - **GPE**: Geopolitical entities (e.g., countries, cities).
- Sentiment analysis was conducted using the TextBlob library to extract sentiment polarity scores.

- The popularity of articles was estimated by counting the number of tweet_ids associated with each headline.

2.3 Feature Engineering

The following features were derived for analysis:

- **NER Features:** Counts of named entities (ORG, PERSON, GPE) in each headline.
- **Article Length:** Number of words in the headline.
- **Sentiment Polarity:** Sentiment scores ranging from -1 (negative) to +1 (positive).
- **Popularity Score:** Number of tweets referencing the article.

2.4 Correlation and Visualization

- A correlation matrix was computed to identify relationships between features such as NER counts, sentiment, article length, and popularity.
- Visualizations included bar charts (comparing NER counts between fake and real news) and scatter plots (popularity vs sentiment).

3. Findings

3.1 Named Entities and Fake/Real News

- Real news articles contained significantly higher counts of named entities (ORG, PERSON, GPE) compared to fake news articles.
- Fake news articles often lacked specific entities, relying instead on vague or sensationalized language.

Named Entity Counts Fake News (Mean) Real News (Mean)

ORG Count	0.57	1.34
PERSON Count	0.32	0.68
GPE Count	0.48	1.12

3.2 Sentiment Analysis

- Sentiment polarity was more extreme for fake news (both highly positive and highly negative), while real news had a more neutral sentiment distribution.

Sentiment Polarity Fake News (Mean) Real News (Mean)

Polarity	0.15	0.05
----------	------	------

3.3 Popularity and Engagement

- Real news articles had higher popularity scores (more tweets) on average than fake news articles. This suggests that real news generates greater public engagement, possibly due to its credibility and informational value.

Popularity Score Fake News (Mean) Real News (Mean)

Tweets Count 13.24 27.89

3.4 Correlation Analysis

The correlation matrix revealed the following insights:

- **Popularity and Named Entities:** There was a positive correlation between the count of named entities (ORG, PERSON, GPE) and the popularity score.
- **Popularity and Sentiment:** Sentiment had a weaker correlation with popularity, suggesting that content features (like entities) are stronger predictors of engagement than sentiment alone.
- **Real/Fake Classification:** NER counts were strong indicators for distinguishing fake from real news.

4. Visualizations

4.1 Correlation Heatmap

The correlation matrix, visualized as a heatmap, highlights relationships between features:

- Strong positive correlation: Named entity counts and popularity.
- Negative correlation: Article length and sentiment with fake news.

4.2 NER Feature Counts

Bar chart comparing NER counts for fake and real news:

- Real news articles consistently exhibited higher counts of organizations, people, and geopolitical entities.

4.3 Popularity vs Sentiment Scatter Plot

Scatter plot showing the relationship between sentiment polarity and popularity:

- Real news articles cluster near neutral sentiment and high popularity.
- Fake news articles scatter more widely, with both highly positive and highly negative sentiments.

5. Insights and Implications

5.1 Named Entities as Indicators

- Real news articles often contain specific, verifiable details (e.g., names of organizations, people, and locations), while fake news lacks such specificity. This makes NER a powerful tool for identifying fake news.

5.2 Sentiment Polarization in Fake News

- Fake news tends to use polarizing language to attract attention. This can be leveraged to detect potential misinformation.

5.3 Popularity Trends

- Real news articles generate more engagement (e.g., tweets), likely due to their credibility. However, fake news with sensationalized language can sometimes achieve similar levels of attention.

5.4 Practical Applications

- **Detection Systems:** Automated systems using NER and sentiment features can identify fake news with reasonable accuracy.
- **Content Analysis:** Media organizations can use these insights to improve engagement with real, high-quality content.

6. Conclusion

This study highlights the importance of named entities in distinguishing fake news from real news and their influence on article popularity. By combining NER, sentiment analysis, and engagement metrics, researchers and practitioners can develop robust tools to combat misinformation and promote credible journalism.

7. Future Work

- Explore additional NER categories (e.g., dates, events) for enhanced fake news detection.
- Analyze full article text (beyond headlines) to capture a broader range of features.
- Investigate the role of visual elements (e.g., images) in news engagement and credibility.

8. References

- SpaCy Documentation: <https://spacy.io>
 - TextBlob Documentation: <https://textblob.readthedocs.io/en/dev/>
 - FakeNewsNet Dataset: <https://github.com/KaiDMML/FakeNewsNet>
-