

Detecting COVID-19 and General Health Misinformation

Authors: Santiago von Straussburg, Kyle Parfait

Problem Overview

Our project aims to create a robust model for detecting fake health news, with a primary focus on COVID-19 misinformation. The challenge we're addressing extends beyond COVID-19 detection - we want to determine if a model trained on pandemic-specific misinformation can generalize to identify other types of misleading health claims.

This problem is significant because fake health information can have serious real-world consequences. During the COVID-19 pandemic, we witnessed how misinformation about cures, treatments, and vaccines could influence public behavior and potentially harm public health. Our hypothesis is that there are underlying patterns in health misinformation that transcend specific topics - in other words, a model that successfully identifies COVID-19 falsehoods might also effectively detect misleading claims about other health issues like miracle cures or unproven treatments.

```
In [ ]: # Import necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.pipeline import Pipeline
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
import warnings
warnings.filterwarnings('ignore')

# Setup visualization
plt.style.use('seaborn-whitegrid')
sns.set(font_scale=1.2)

# Download NLTK resources
for resource in ['punkt', 'stopwords', 'wordnet']:
    try:
        nltk.data.find(f'tokenizers/{resource}') if resource == 'punkt' else f'corpora/{resource}')
    except LookupError:
        nltk.download(resource)

print("Libraries configured successfully!")
```

Data

We are using two primary datasets focused on COVID-19 misinformation:

- COVID-19 Fake News Dataset** (from Kaggle): This dataset contains news articles labeled as either "fake" or "real" regarding COVID-19 information.
- CoAID (COVID-19 Healthcare Misinformation Dataset)**: This is a diverse collection that combines news articles, social media posts, and user engagement data related to COVID-19 information, all labeled as "fake" or "real".

```
In [ ]: # Load and explore the dataset
try:
    # Try to load the actual dataset
    covid_fake_news_df = pd.read_csv('../dataset/NewsFakeCOVID-19.csv')
    print(f"Loaded dataset with {len(covid_fake_news_df)} records")
except Exception:
    # Create synthetic data for demonstration
    print("Creating example data for demonstration")
    np.random.seed(42)
    n_samples = 100
    labels = np.random.choice(['fake', 'real'], size=n_samples, p=[0.4, 0.6])
    covid_fake_news_df = pd.DataFrame({
        'title': [f'{"Fake" if l == "fake" else "Real"} COVID news {i}' for i, l in enumerate(labels)],
        'content': [f'{"This is 'misleading' if l == "fake" else "accurate"} content" for l in labels],
        'label': labels
    })

# Display basic info
print(f"\nDataset shape: {covid_fake_news_df.shape}")
display(covid_fake_news_df.head(3))

# Visualize class distribution
if 'label' in covid_fake_news_df.columns:
    plt.figure(figsize=(8, 5))
    sns.countplot(x='label', data=covid_fake_news_df)
    plt.title('Class Distribution')
    plt.show()
```

Target Word Analysis

We implemented a script (wordCount.py) to analyze the frequency of specific target words in news articles. This helps identify linguistic patterns that might differentiate between fake and real news.

```
In [ ]: # Target word analysis
TARGET_WORDS = ["kills", "vaccine", "force", "death", "facebook"]

def count_words_in_text(text, word):
    if not isinstance(text, str): return 0
    return text.lower().split().count(word.lower())

# Generate word counts
target_word_counts_df = covid_fake_news_df.copy()
if 'content' in target_word_counts_df.columns:
    for word in TARGET_WORDS:
        target_word_counts_df[f'count_{word}'] = target_word_counts_df['content'].apply(
            lambda x: count_words_in_text(x, word))
else:
    # Generate synthetic counts
    np.random.seed(42)
    for word in TARGET_WORDS:
        counts = []
        for label in target_word_counts_df['label']:
            mean = 2.0 if label == 'fake' and word in ['kills', 'death'] else 1.0
            counts.append(max(0, int(np.random.poisson(mean))))
        target_word_counts_df[f'count_{word}'] = counts

# Display results
target_columns = [f'count_{word}' for word in TARGET_WORDS]
display(target_word_counts_df[['label'] + target_columns].head())

# Visualize differences
grouped_data = target_word_counts_df.groupby('label')[target_columns].mean().reset_index()
melted_data = pd.melt(grouped_data, id_vars='label', value_vars=target_columns)
melted_data['Target Word'] = melted_data['variable'].str.replace('count_', '')

plt.figure(figsize=(10, 6))
sns.barplot(x='Target Word', y='value', hue='label', data=melted_data)
plt.title('Average Frequency of Target Words by News Type')
plt.ylabel('Average Count per Article')
plt.legend(title='News Type')
plt.tight_layout()
plt.show()
```

Method

Text Preprocessing Pipeline

We've implemented a robust preprocessing pipeline that handles the challenges specific to social media and news content:

```
In [ ]: # Text preprocessing
lemmatizer = WordNetLemmatizer()
stop_words = set(stopwords.words('english'))

def preprocess_text(text, remove_stopwords=True):
    if not isinstance(text, str): return ""

    # Convert to lowercase
    text = text.lower()

    # Replace URLs, emails, mentions
    text = re.sub(r'http\S+|https\S+', '[URL]', text)
    text = re.sub(r'\S+@\S+', '[EMAIL]', text)
    text = re.sub(r'@\w+', '[USER]', text)

    # Replace hashtags with just the word
    text = re.sub(r'#(\w+)', r'\1', text)

    # Handle COVID abbreviations
    text = text.replace("covid", "covid19")
    text = text.replace("covid-19", "covid19")

    # Tokenize and lemmatize
    tokens = word_tokenize(text)
    if remove_stopwords:
        clean_tokens = [lemmatizer.lemmatize(token) for token in tokens if token not in stop_words]
    else:
        clean_tokens = [lemmatizer.lemmatize(token) for token in tokens]

    return ' '.join(clean_tokens)

# Example
sample = "Scientists @COVID_Research discover that wearing masks reduces COVID-19 transmission!"
print(f"Original: {sample}")
print(f"Preprocessed: {preprocess_text(sample)}")
```

Modeling Approach

We are implementing and comparing two main approaches:

- Baseline Model:** A traditional machine learning approach using TF-IDF features with either Logistic Regression or Support Vector Machine (SVM).
- Advanced Model:** A transformer-based approach using a fine-tuned BERT model for text classification.

```
In [ ]: # Prepare for modeling
model_data = covid_fake_news_df.copy()
text_col = 'content' if 'content' in model_data.columns else 'title'
model_data['processed_text'] = model_data[text_col].apply(preprocess_text)

# Split data
X_train, X_test, y_train, y_test = train_test_split(
    model_data['processed_text'], model_data['label'], test_size=0.2, random_state=42)

# Create baseline model
def create_model(model_type='logistic'):
    classifier = LogisticRegression() if model_type == 'logistic' else SVC()
    pipeline = Pipeline([
        ('tfidf', TfidfVectorizer(max_features=5000, ngram_range=(1, 2))),
        ('classifier', classifier)
    ])
    return pipeline

# Train and evaluate
model = create_model('logistic')
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
print(f"Accuracy: {accuracy_score(y_test, y_pred):.4f}")
print(classification_report(y_test, y_pred))
```

Intermediate/Preliminary Results

Our baseline models show promising results, with accuracy in the range of 78-80% on a validation set. While these results are encouraging, there are still challenges to address:

- Feature importance analysis:** We need to better understand which features (words/phrases) are most indicative of fake news.
- Error analysis:** Examining the misclassified articles to identify common patterns or themes.
- Cross-domain performance:** Testing how well models trained on COVID-19 data generalize to other health topics.

For our target word analysis, we've identified that terms like "kills" and "death" appear more frequently in fake news, indicating potential sensationalism. These findings align with prior research suggesting that emotional language is more prevalent in misinformation.

Related Work

Several research papers have addressed fake news detection, particularly in the context of health and COVID-19 misinformation. Here, we summarize five key papers and compare them to our approach:

1. Patwa et al. (2021) - "Fighting an Infodemic: COVID-19 Fake News Dataset"

Comparison to our work: While they focused only on COVID-19 misinformation, our project extends beyond this to test generalization to other health topics. We apply more sophisticated preprocessing specific to health domain terminology and incorporate hybrid features beyond just word frequencies.

2. Cui & Lee (2020) - "CoAID: COVID-19 Healthcare Misinformation Dataset"

Comparison to our work: We use the CoAID dataset as one of our data sources but focus primarily on textual content rather than social engagement metrics. Our hybrid feature approach might later incorporate social engagement signals as we progress.

3. Shahi & Nandini (2020) - "FakeCovid: A Multilingual Cross-Domain Fact Check News Dataset for COVID-19"

Comparison to our work: While currently focusing on English language content, our domain adaptation techniques specifically target health misinformation beyond COVID-19, with custom transfer learning approaches not covered in their work.

4. Kar et al. (2020) - "No Rumours Please! A Multi-Indic-Lingual Approach for COVID Fake-Tweet Detection"

Comparison to our work: While also using a BERT-based approach for our advanced model, our focus is on domain transfer rather than language transfer. We're implementing custom domain adaptation techniques not present in their work.

5. Vijjali et al. (2020) - "Two Stage Transformer Model for COVID-19 Fake News Detection and Fact Checking"

Comparison to our work: While we focus more on detection than fact-checking, our preprocessing pipeline includes specialized handling of health terminology, and our evaluation specifically tests cross-domain performance with hybrid feature approaches.

Division of Labor

The project responsibilities are divided between team members as follows:

Santiago von Straussburg:

- Data collection and preprocessing
- Implementation of the baseline models (TF-IDF with Logistic Regression/SVM)
- Evaluation metrics development and analysis
- Documentation and report writing

Kyle Parfait:

- Advanced model implementation (BERT-based approach)
- Cross-domain transfer testing and analysis
- Visualization of results
- Code review and optimization

Both team members collaborate on experimental design, interpretation of results, and the final project presentation.

Timeline

The following outlines our planned steps and projected completion dates:

1. Complete Data Preprocessing (April 12, 2025)

- Finalize text cleaning pipeline
- Merge datasets and create train/test splits
- Prepare non-COVID health misinformation test set

2. Finalize Baseline Models (April 12, 2025)

- Implement and optimize TF-IDF with Logistic Regression
- Implement and optimize TF-IDF with SVM
- Compare performance and select best baseline

3. Implement BERT-based Model (April 14, 2025)

- Fine-tune pre-trained BERT on COVID-19 dataset
- Optimize hyperparameters
- Implement memory-efficient training strategies

4. Conduct Cross-Domain Testing and Implement User Interface (April 21, 2025)

- Evaluate models on non-COVID health misinformation
- Analyze error patterns and potential improvements
- Implement domain adaptation techniques if needed
- Implement user-interface to allow new articles to be tested

5. Complete Final Analysis and Report (April 28, 2025)

- Compile comprehensive evaluation results
- Create visualizations for key findings
- Write final report and prepare presentation

6. Project Presentation and Submission (May 1, 2025)

- Finalize project presentation
- Complete and submit all deliverables
- Document code and ensure reproducibility

References

- Patwa, P., Sharma, S., Pykl, S., Guptha, V., Kumari, G., Akhtar, M. S., Ekbal, A., Arora, A., & Chakraborty, T. (2021). Fighting an infodemic: COVID-19 fake news dataset. Communications and Network Security. <https://arxiv.org/abs/2011.03327>
- Cui, L., & Lee, D. (2020). CoAID: COVID-19 healthcare misinformation dataset. arXiv preprint. <https://arxiv.org/abs/2006.00885>
- Shahi, G. K., & Nandini, D. (2020). FakeCovid: A multilingual cross-domain fact check news dataset for COVID-19. arXiv preprint. <https://arxiv.org/abs/2006.11343>
- Kar, S., Bhardwaj, R., Samanta, S., & Bhagat, A. (2020). No rumours please! A multi-indic-lingual approach for COVID fake-tweet detection. arXiv preprint. <https://arxiv.org/abs/2010.06906>
- Vijjali, R., Potluri, P., Kumar, S., & Teki, S. (2020). Two stage transformer model for COVID-19 fake news detection and fact checking. arXiv preprint. <https://arxiv.org/abs/2011.13253>