# Data Collection and Preprocessing Phase

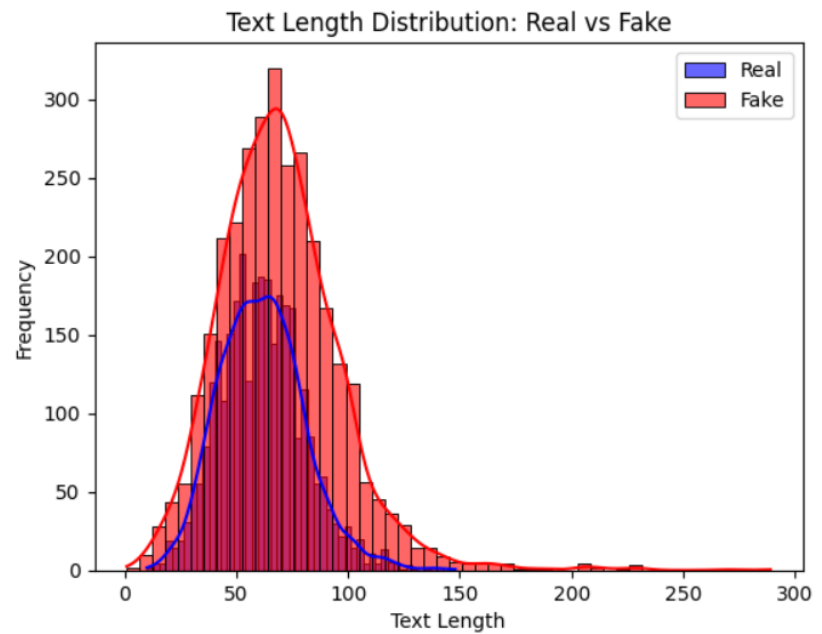| | |
|---|---|
| Date | 23 September 2024 |
| Team ID | LTVIP2024TMID25030 |
| Project Title | FAKE NEWS ANALYSIS IN SOCIAL MEDIA |
| Maximum Marks | 6 Marks |

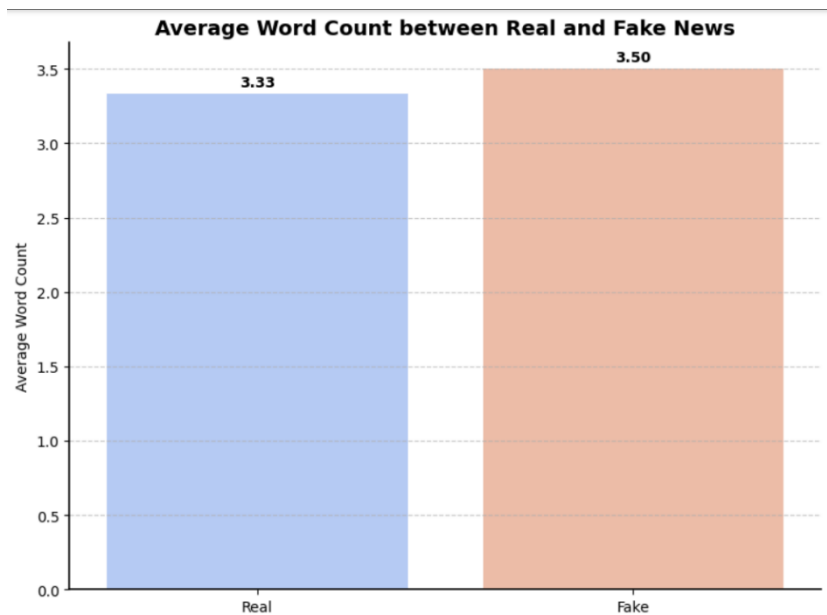**Data Exploration and Preprocessing Template**

Dataset variables will be statistically analysed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modelling, and forming a strong foundation for insights and predictions.
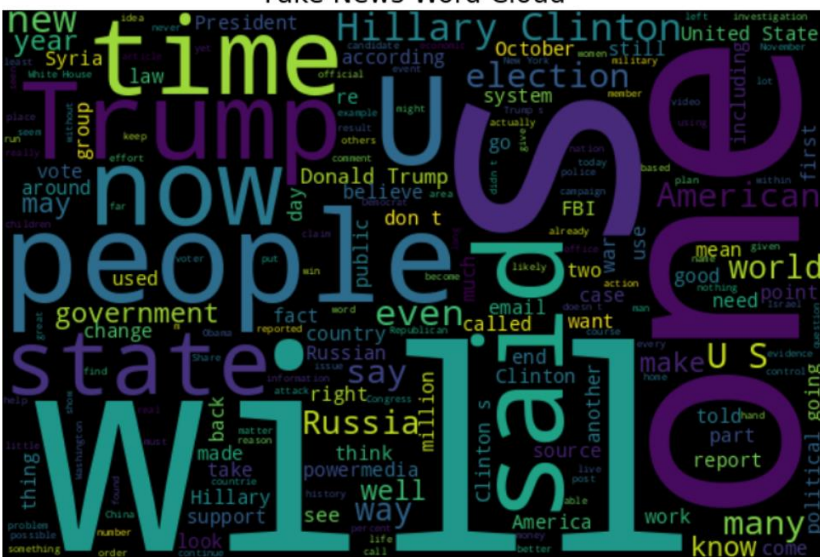
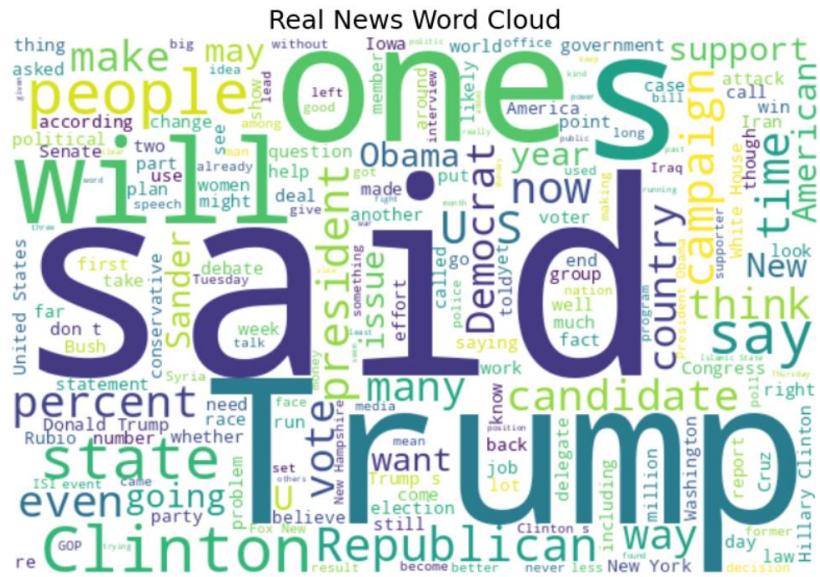| Section | Description |
|---|---|
| Data Overview | Dimension:<br>9285 Rows X 4 Columns<br>Descriptive statistics:<br> |

| | |
|---|---|
| Univariate Analysis | Text Length Distribution: Real vs Fake |
| Bivariate Analysis | Average Word Count between Real and Fake News |

| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies | - |
| **Data Preprocessing Code Screenshots** | |
| Loading Data | ```
# Load the dataset
df=pd.read_csv("/content/drive/MyDrive/Fake News.csv",usecols=['title','text','label'],encoding='latin1',on_bad_lines='skip')
df.head()
``` |

| | |
|---|---|
| | |
| Handling Missing Data | ```python
# Dropping unnecessary columns (modify if necessary)
df.drop(columns=['Unnamed: 0'], inplace=True,errors='ignore')

# Checking for missing values
print(df.isnull().sum())
``` |
| Data Transformation | ```python
import re
import string

# Text cleaning function
def clean_text(text):
    # Check if text is a string before applying lower()
    if isinstance(text, str):
        text = text.lower()  # Convert to lowercase
        text = re.sub('\[.*?\]', '', text)  # Remove text in square brackets
        text = re.sub('https?://\S+|www\.\S+', '', text)  # Remove links
        text = re.sub('<.*?>+', '', text)  # Remove HTML tags
        text = re.sub('[%s]' % re.escape(string.punctuation), '', text)  # Remove punctuation
        text = re.sub('\n', '', text)  # Remove newline characters
        text = re.sub('\w*\d\w*', '', text)  # Remove words containing numbers
        return text
    else:
        # Handle non-string values (e.g., return empty string or NaN)
        return ''  # or return float('nan')

# Apply the cleaning function to the 'text' column
df['cleaned_text'] = df['text'].apply(clean_text)

# Display the cleaned text
print(df[['text', 'cleaned_text']].head())
``` |
| Feature Engineering | Attached the code in the final submission |
| Save Processed Data | - |