

# Coursework Assignment: Text classification

## I. Introduction

As I'm deeply interested in the intersection of technology and media, I chose to explore the accuracy of news categorization in digital platforms. The proliferation of digital news media has led to an immense volume of articles categorized into various topics. However, the accuracy of these categorizations can significantly impact information retrieval and reader experience. This project aims to validate and potentially refine the categorization of news articles in the "Daily Google News" dataset. By applying text classification techniques, the project seeks to evaluate the precision of current categorizations and explore the possibility of more nuanced sub-categories, enhancing the dataset's utility for both academic and practical applications.

### 1. Problem Area

**Context and Relevance:** The digital age has witnessed an exponential increase in online news content, making it challenging for readers to navigate and extract value effectively. This abundance poses a significant challenge: how to categorize and analyze vast amounts of news data efficiently. The need for advanced automated methods to manage and understand this content is more pressing than ever. This is particularly crucial in an era where information overload is a common concern, and the rapid identification of relevant news is vital for both individuals and organizations.

**Specific Challenge:** The specific challenge addressed in this coursework lies in the realm of automated text classification, focusing on news titles. News titles are succinct, yet they carry the essence of the articles they represent. However, their brevity and often sensational nature can make classification and analysis non-trivial. The task is to develop a system capable of categorizing these titles into meaningful groups (such as politics, sports, business, etc.) and performing further analysis like trend identification and sentiment analysis. This approach is not only innovative but also highly relevant in today's fast-paced information cycle, where quick and accurate categorization of news can lead to better information management and decision-making.

**Industry and Societal Impact:** The implications of this project extend beyond academic interest; they have real-world applications in various sectors. For media outlets and news aggregators, efficient categorization can enhance user experience and content delivery. For businesses and analysts, trend analysis can inform market strategies and public sentiment assessment. Additionally, in the societal context, such a system can aid in combating the spread of misinformation by quickly categorizing and flagging potential fake news or biased reporting.

**Building on Previous Work:** This project is inspired by existing research in Natural Language Processing (NLP) and machine learning, particularly in text classification and

sentiment analysis. Previous studies have laid the groundwork by exploring various models and techniques, including but not limited to Support Vector Machines, Naïve Bayes classifiers, and neural networks. However, this project aims to innovate by focusing on the concise and often nuanced nature of news titles, a less explored domain compared to full-article analysis.

**Novel Contribution:** The novelty of this project lies in its targeted focus on news titles using advanced NLP techniques. While previous work in the field has primarily concentrated on longer texts, this project recognizes the unique challenges and opportunities presented by the concise nature of titles. By employing and potentially enhancing state-of-the-art NLP methodologies, this project aims to contribute to the field by offering a specialized solution tailored to the peculiarities of news title classification and analysis.

## 2. Objectives

Develop a sophisticated text classification framework tailored for the "Daily Google News" dataset. This framework is not just focused on categorizing news titles accurately but extends to deeper analyses such as trend identification, sentiment evaluation, and detailed examination of language styles.

### Detailed Objectives:

1. **Comprehensive Categorization:** The framework will employ advanced classification algorithms to segment news titles into clearly defined categories like politics, sports, and business. This step forms the foundation of the analysis, capturing the wide-ranging themes and diversity within the dataset.
2. **Trend Analysis:** After initial categorization, the system will examine temporal trends. This includes tracking how certain news categories fluctuate over time, offering insights into changing patterns and priorities in news reporting. This analysis will help to understand how public interest or media focus shifts in different periods.
3. **Sentiment Analysis:** The framework will apply sentiment analysis techniques to evaluate the emotional undertones of news titles. This objective aims to uncover potential biases or editorial slants in the way news is presented, contributing to a deeper understanding of media psychology and content framing.
4. **Language and Style Analysis:** An in-depth linguistic analysis will be conducted to examine the stylistic elements and structural aspects of news titles. This will provide insights into the editorial choices and language use across various news platforms, emphasizing the role of language in shaping news appeal and perception.

### Significance and Impact:

- This project will significantly enhance the precision and depth of news title analysis, moving beyond simple categorization to a more nuanced and comprehensive understanding of news content.
- By offering a multi-faceted perspective of the news data, the framework will highlight crucial aspects like evolving trends, embedded sentiments, and distinctive editorial styles, crucial in today's digital news landscape.

- The insights gained from this analysis will be immensely beneficial for digital news platforms, improving content curation and user experience, and contributing to more responsible and accurate news dissemination.

### **Novelty and Methodological Approach:**

- The framework will integrate traditional Natural Language Processing (NLP) techniques with advanced artificial intelligence models to address the unique challenges posed by the concise and varied nature of news titles.
- Innovative approaches in unsupervised learning and semantic analysis will be explored to uncover hidden sub-categories and deeper meanings within the dataset, pushing the boundaries of conventional text classification.
- A diverse set of analytical tools and methodologies will be employed to offer a holistic view of the dataset, showcasing not only technical prowess but also a profound understanding of media studies and digital communication.

## **3. Dataset**

In this project, I worked with the '**Daily Google News**' dataset, sourced from [Kaggle](#). This dataset comprises a comprehensive collection of news articles aggregated over several months in 2023. It includes four distinct CSV files, each corresponding to a specific month. The dataset's structure and content are outlined as follows:

### **1. Monthly Files:**

- **September 2023 (2023\_9.csv):** Contains 44,834 unique news articles.
- **October 2023 (2023\_10.csv):** Encompasses 46,261 unique news articles.
- **November 2023 (2023\_11.csv):** Includes 36,178 unique news articles.
- **December 2023 (2023\_12.csv):** Features 15,735 unique news articles.

### **2. Data Columns:** Each CSV file consists of five key columns:

- **Title:** The headline of the news article.
- **Publisher:** The source or publication that released the news article.
- **DateTime:** The date and time when the article was published on Google News.
- **Link:** A URL leading to the full news article, enabling further exploration and content analysis.
- **Category:** The news category, as defined by Google News, covering major domains such as Business, Entertainment, Headlines, Health, Science, Sports, Technology, and WorldWide.

### **3. Diversity and Scope:** The dataset presents a diverse range of topics and publishers, offering a broad perspective on global and local news events. The categorization by Google News into various domains allows for targeted analysis in specific fields of interest.

### **4. Temporal Coverage:** Spanning from September to December 2023, the dataset provides a snapshot of the news landscape over these months. This temporal range is valuable for identifying trends, seasonal variances, and shifts in media focus.

### **5. Potential for Deep Diving:** The inclusion of direct links to the original articles opens opportunities for extended analysis, including content scraping for more detailed textual data, sentiment extraction, and deeper linguistic studies.

6. **Significance for Analysis:** The dataset's comprehensive nature makes it a fertile ground for our text classification project. The richness in topics, coupled with the temporal breadth, allows for an extensive examination of news trends, categorization efficiency, and style analysis.

## 4. Evaluation Methodology

The evaluation methodology for our text classification project is focused on rigorously assessing the model's performance in categorizing news titles and capturing various analytical dimensions such as trends, sentiments, and linguistic styles. Here's an overview of our comprehensive approach:

### 1. Selection of Metrics:

- **Accuracy:** Measures the overall effectiveness of the classification model.
- **Precision and Recall:** Crucial for imbalanced datasets to understand the model's performance for underrepresented categories.
- **F1-Score:** Balances precision and recall, providing a comprehensive view of model performance.
- **Confusion Matrix:** Offers insights into the model's performance across different categories, highlighting potential biases or misclassification tendencies.

### 2. Objective-Specific Metrics:

- **Trend Analysis:** Evaluates the model's ability to track category prevalence over time.
- **Sentiment Analysis:** Measures the accuracy of sentiment categorization.
- **Language and Style Analysis:** Assesses how well the model captures linguistic features and styles.

### 3. Validation Strategy:

- **Cross-Validation:** Utilizes K-fold cross-validation for consistent performance across different dataset subsets.
- **Train-Test Split:** Separates a portion of the dataset for unbiased final model testing.

### 4. Benchmarking:

- **Comparative Analysis:** The model's performance will be compared against benchmarks established in the paper "[Automatic Semantic Categorization of News Headlines using Ensemble Machine Learning: A Comparative Study](#)" (IJACSA, Vol. 10, No. 11, 2019). This study achieved an accuracy of 90.12% and recall of 90% with Multinomial Naïve Bayes, setting a high standard for our model.

### 5. Error Analysis:

- Conducts a thorough examination of instances where the model underperforms to identify limitations and improvement areas.

## II. Implementation

## 5. Preprocessing

```
In [1]: # ! pip install pandas numpy nltk scikit-learn seaborn textstat textblob
```

```
In [2]: # Import necessary libraries
import pandas as pd
import re
import nltk
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.probability import FreqDist
from nltk.util import ngrams
from nltk.stem import WordNetLemmatizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import classification_report, accuracy_score

from textblob import TextBlob

# Download necessary NLTK data
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')
nltk.download('maxent_ne_chunker')
nltk.download('words')

# Load dataset
csv_files = ['./2023_9.csv', './2023_10.csv', './2023_11.csv', './2023_12.csv']
dataframes = [pd.read_csv(file) for file in csv_files]
data = pd.concat(dataframes, ignore_index=True)

data.dropna(subset=['Title'])

# Print the column names of the DataFrame
print(data.columns)

# Step 1: Data Cleaning and Normalization
def clean_and_normalize(text):
    # Check if the text is a string
    if isinstance(text, str):
        text = re.sub(r'<.*?>', '', text) # Remove HTML tags
        text = re.sub(r'^\w\s', '', text) # Remove punctuation
        text = text.lower() # Convert to lowercase
        return text
    else:
        return "" # Return empty string for non-string inputs

data['title'] = data['Title'].apply(clean_and_normalize)

# Step 2: Text Tokenization
def tokenize_text(text):
    return word_tokenize(text)

data['tokens'] = data['title'].apply(tokenize_text)

# Step 3: Stop Word Removal
stop_words = set(stopwords.words('english'))
def remove_stopwords(tokens):
```

```
    return [word for word in tokens if word not in stop_words]

data['filtered_tokens'] = data['tokens'].apply(remove_stopwords)

# Display the DataFrame with preprocessing steps applied

data.head( )
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]   date!
[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package maxent_ne_chunker is already up-to-date!
[nltk_data] Downloading package words to
[nltk_data]   /Users/heyonggang/nltk_data...
[nltk_data]   Package words is already up-to-date!
Index(['Title', 'Publisher', 'DateTime', 'Link', 'Category'], dtype='object')
```

Out [2]:

|   | Title   | Publisher              | DateTime             | Link  |
|---|---|------------------------|----------------------|---|
| 0 | Chainlink (LINK)alters, Hedera (HBAR) Wobble...   | Analytics Insight      | 2023-08-30T06:54:49Z | <a href="https://news.google.com/articles/CBMibGh0dHBzO...">https://news.google.com/articles/CBMibGh0dHBzO...</a> |
| 1 | Funds punished for owning too few Nvidia share... | ZAWYA                  | 2023-08-30T07:15:59Z | <a href="https://news.google.com/articles/CBMigwFodHRwc...">https://news.google.com/articles/CBMigwFodHRwc...</a> |
| 2 | Crude oil prices stalled as hedge funds sold: ... | ZAWYA                  | 2023-08-30T07:31:31Z | <a href="https://news.google.com/articles/CBMibGh0dHBzO...">https://news.google.com/articles/CBMibGh0dHBzO...</a> |
| 3 | Grayscale's Bitcoin Win Is Still Only Half the... | Bloomberg              | 2023-08-30T10:38:40Z | <a href="https://news.google.com/articles/CBMib2h0dHBzO...">https://news.google.com/articles/CBMib2h0dHBzO...</a> |
| 4 | I'm a Home Shopping Editor, and These Are the ... | Better Homes & Gardens | 2023-08-30T11:00:00Z | <a href="https://news.google.com/articles/CBMiPWh0dHBzO...">https://news.google.com/articles/CBMiPWh0dHBzO...</a> |

## 5.1 News Title Categorization

In [3]:

```
# Feature Extraction using TF-IDF
tfidf_vectorizer = TfidfVectorizer()
X_features = tfidf_vectorizer.fit_transform(data['title'])

# Naive Bayes classifier for text classification
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB

# Splitting the dataset
X_train, X_test, y_train, y_test = train_test_split(X_features, data['Category'])

# Model Training
model = MultinomialNB()
model.fit(X_train, y_train)

# Evaluate the model's performance on the test set using accuracy and other rel
from sklearn.metrics import classification_report

# Predictions
y_pred = model.predict(X_test)

# Accuracy and Classification Report
```

```
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
```

Accuracy: 0.7374792186201163

Classification Report:

|               | precision | recall | f1-score | support |
|---------------|-----------|--------|----------|---------|
| Business      | 0.86      | 0.69   | 0.76     | 5197    |
| Entertainment | 0.90      | 0.78   | 0.83     | 5217    |
| Headlines     | 0.71      | 0.53   | 0.61     | 5468    |
| Health        | 0.89      | 0.75   | 0.82     | 3452    |
| Science       | 0.94      | 0.72   | 0.81     | 3297    |
| Sports        | 0.50      | 0.97   | 0.66     | 5825    |
| Technology    | 0.90      | 0.78   | 0.83     | 5020    |
| Worldwide     | 0.70      | 0.67   | 0.69     | 5020    |
| accuracy      |           |        | 0.74     | 38496   |
| macro avg     | 0.80      | 0.74   | 0.75     | 38496   |
| weighted avg  | 0.78      | 0.74   | 0.74     | 38496   |

The accuracy of about 74% indicates a reasonable performance for this model. However, there are some categories where the recall is comparatively low, such as in 'Headlines', 'Science', 'Sports' and 'Worldwide' suggesting that the model is not identifying these categories as effectively as others.

Word2Vec can capture semantic relationships between words better than traditional bag-of-words models Which can produce more accurate results.

```
In [4]: # Install Gensim
!pip install gensim
# Import the necessary libraries:
import gensim
from gensim.models import Word2Vec

# Train the Word2Vec Model
data1 = data.copy()
model = Word2Vec(sentences=data1['filtered_tokens'], vector_size=100, window=5,

# Save the trained model for future use:
model.save("word2vec.model")

# Load saved model
model = Word2Vec.load("word2vec.model")

# Use the model to obtain word vectors
vector = model.wv['example'] # Get vector for a word
similar_words = model.wv.most_similar('example') # Find similar words
```

Requirement already satisfied: gensim in /usr/local/lib/python3.11/site-packages (4.3.2)

Requirement already satisfied: numpy>=1.18.5 in /usr/local/lib/python3.11/site-packages (from gensim) (1.26.3)

Requirement already satisfied: scipy>=1.7.0 in /usr/local/lib/python3.11/site-packages (from gensim) (1.11.4)

Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.11/site-packages (from gensim) (6.4.0)

Convert Sentences to Vectors



```
In [5]: import numpy as np

# Function to convert a sentence to a vector by averaging the vectors of the words
def sentence_to_vector(sentence, model):
    # Filter words in sentence if they are in the model's vocabulary
    words = [word for word in sentence if word in model.wv.key_to_index]
    if len(words) == 0:
        return np.zeros(model.vector_size) # Return zero vector if no words are found
    word_vectors = [model.wv[word] for word in words]
    sentence_vector = np.mean(word_vectors, axis=0)
    return sentence_vector

# Convert each sentence in the dataset to a vector
X_vectors = np.array([sentence_to_vector(sentence, model) for sentence in data1])
```

## Classifier Training and Evaluation

```
In [6]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X_vectors, data1['Category'])

# Train a classifier, Random Forest
classifier = RandomForestClassifier(n_estimators=100, random_state=42)
classifier.fit(X_train, y_train)

# Make predictions on the test set
y_pred = classifier.predict(X_test)

# Evaluate the model's performance
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))
```

Accuracy: 0.6448202410640067

Classification Report:

|               | precision | recall | f1-score | support |
|---------------|-----------|--------|----------|---------|
| Business      | 0.72      | 0.59   | 0.65     | 5197    |
| Entertainment | 0.75      | 0.66   | 0.70     | 5217    |
| Headlines     | 0.55      | 0.42   | 0.48     | 5468    |
| Health        | 0.84      | 0.70   | 0.77     | 3452    |
| Science       | 0.87      | 0.70   | 0.78     | 3297    |
| Sports        | 0.46      | 0.88   | 0.61     | 5825    |
| Technology    | 0.82      | 0.69   | 0.75     | 5020    |
| Worldwide     | 0.59      | 0.53   | 0.56     | 5020    |
| accuracy      |           |        | 0.64     | 38496   |
| macro avg     | 0.70      | 0.65   | 0.66     | 38496   |
| weighted avg  | 0.68      | 0.64   | 0.65     | 38496   |

## 5.2 News Title Trend Analysis

```
In [7]: import pandas as pd
import matplotlib.pyplot as plt
import datetime
```

```

# Reload the data with 'DateTime' column
csv_files = ['./2023_9.csv', './2023_10.csv', './2023_11.csv', './2023_12.csv']
dataframes = [pd.read_csv(file) for file in csv_files]
data = pd.concat(dataframes, ignore_index=True)

# Ensure 'DateTime' column is present
print("Columns in the dataset: ", data.columns)

# Convert 'DateTime' to datetime object and set as index
data['DateTime'] = pd.to_datetime(data['DateTime'], format='%Y-%m-%dT%H:%M:%SZ')
data.set_index('DateTime', inplace=True)

# Aggregate Data by Day and Category
daily_trends = data.groupby([pd.Grouper(freq='D'), 'Category']).size().unstack()

# Convert daily counts to cumulative count for a stacked area chart
stacked_data = daily_trends.cumsum(axis=1)

# Creating a Stacked Area Chart
plt.figure(figsize=(15, 8))
stacked_data.plot(kind='area', stacked=True, figsize=(15, 8))

# Adding title and labels to the chart
plt.title('Cumulative Daily News Category Trends (September – December 2023)')
plt.xlabel('Date')
plt.ylabel('Cumulative Number of Articles')

# Setting x-axis limits from 2023-09-01 to 2023-12-31
plt.xlim(datetime.date(2023, 9, 1), datetime.date(2023, 12, 31))

# Rotating the x-axis labels for better readability
plt.xticks(rotation=45)

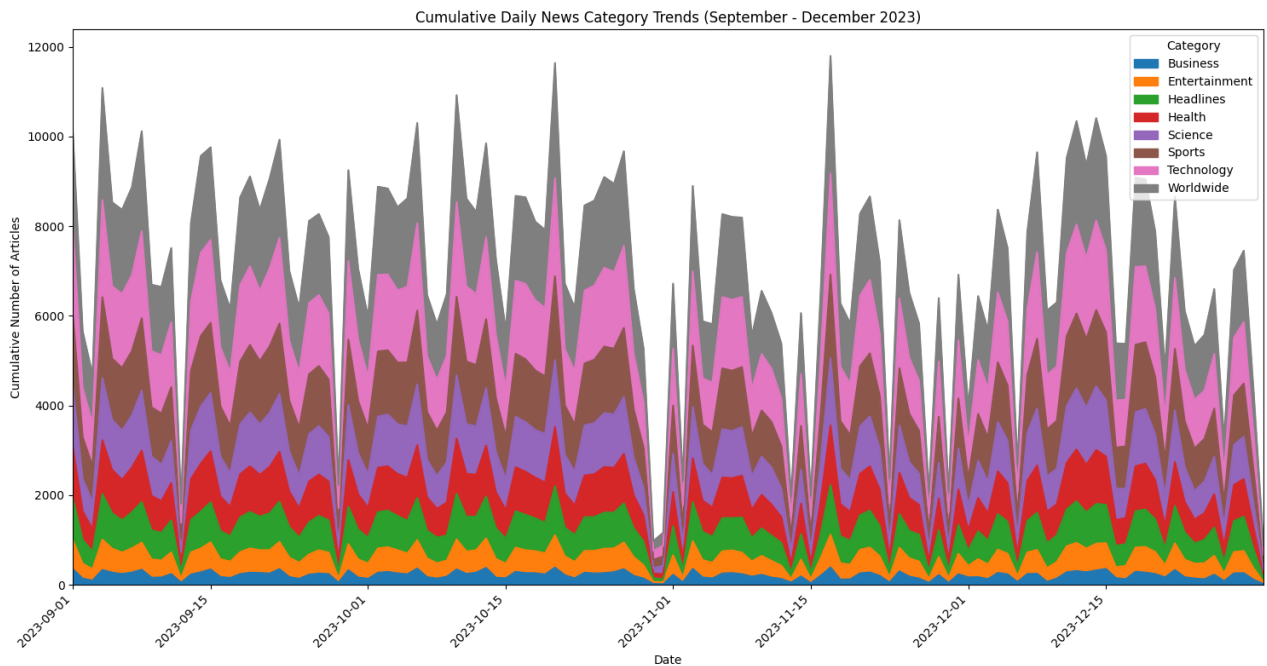
# Applying tight layout for neatness
plt.tight_layout()

# Display the plot
plt.show()

```

Columns in the dataset: Index(['Title', 'Publisher', 'DateTime', 'Link', 'Category'], dtype='object')

<Figure size 1500x800 with 0 Axes>



### 5.3 News Title Sentiment Analysis

1. Preprocessing Since sentiment analysis often depends on the context provided by stop words, I'll avoid removing them in this case. However, I will still perform basic preprocessing like lowercasing and removing special characters.

```
In [8]: def preprocess_for_sentiment(text):
# Check if the text is a string
if isinstance(text, str):
    text = text.lower() # Lowercasing
    text = re.sub(r'^\w\s', '', text) # Removing special characters
    return text
else:
    return "" # Return empty string for non-string inputs

data['title_for_sentiment'] = data['Title'].apply(preprocess_for_sentiment)
```

2. Sentiment Analysis For sentiment analysis, I'll use a pre-trained model from a library like TextBlob or VADER, which are specifically designed for sentiment analysis and are easy to implement.

```
In [9]: from textblob import TextBlob

def analyze_sentiment(text):
    return TextBlob(text).sentiment.polarity

data['sentiment'] = data['title_for_sentiment'].apply(analyze_sentiment)
```

3. Categorizing Sentiment To categorize the sentiment, I'll classify each score into 'positive', 'neutral', or 'negative' based on its value.

```
In [10]: def categorize_sentiment(score):
if score > 0.1:
```

```

        return "Positive"
    elif score < -0.1:
        return "Negative"
    else:
        return "Neutral"

data['sentiment_category'] = data['sentiment'].apply(categorize_sentiment)

```

4. Analysis With the sentiment categories in place, I can analyze the distribution of sentiments across different news categories or over time. To visualize the results, I can use libraries like matplotlib or seaborn to create graphs showing sentiment distribution.

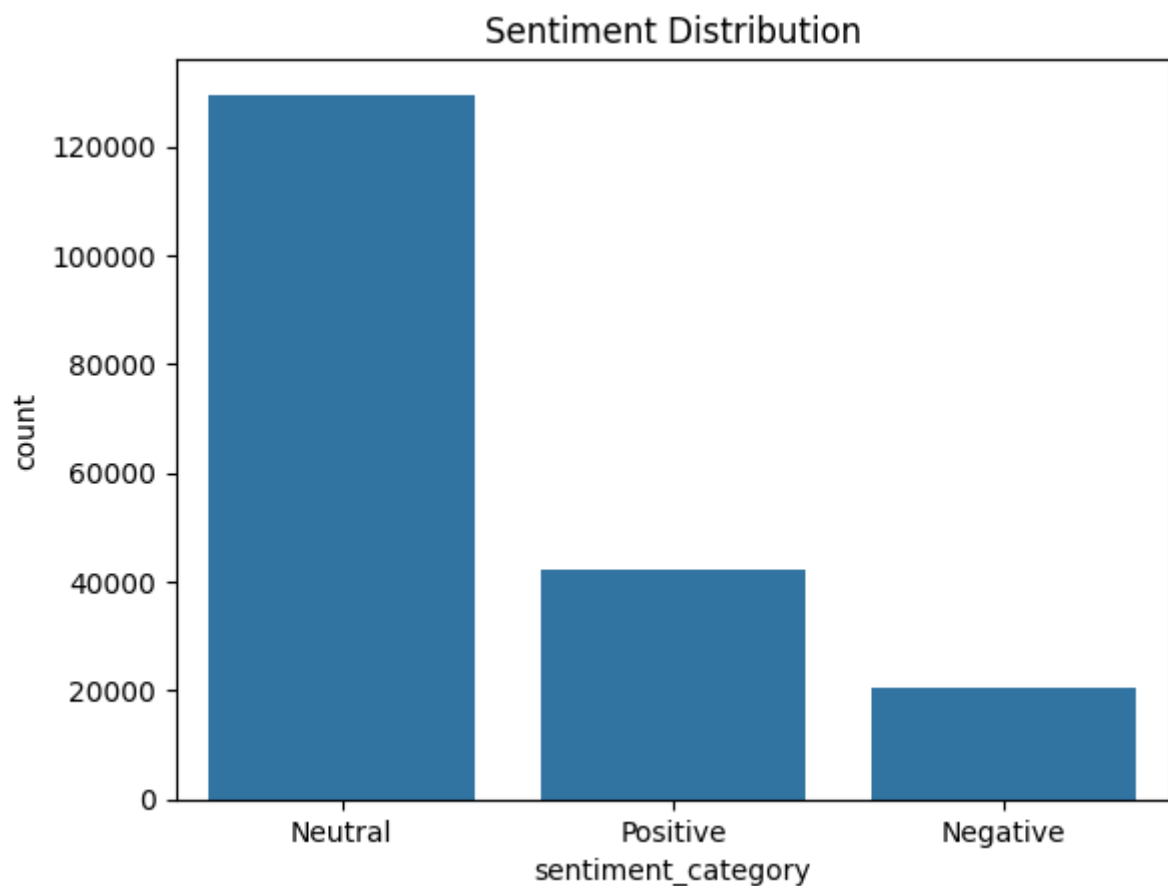
```

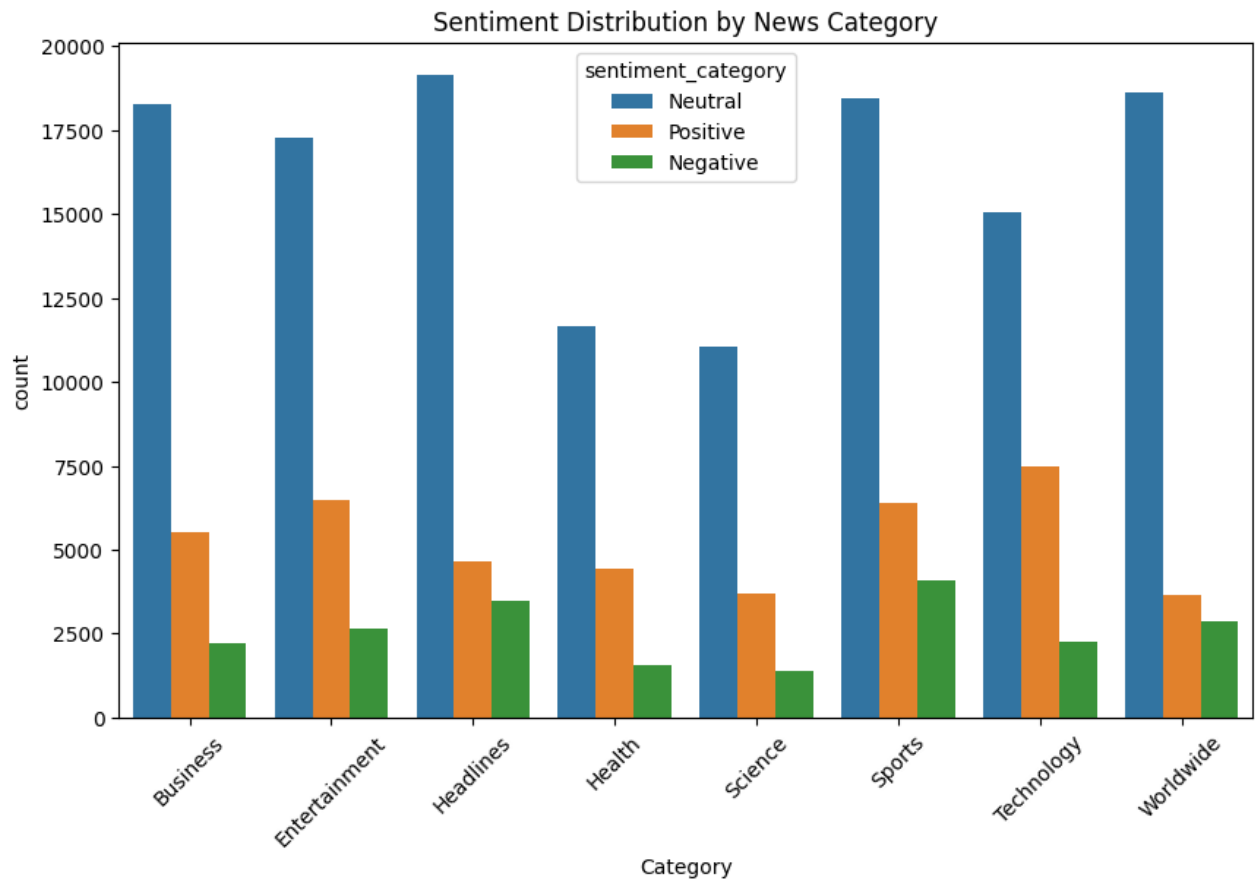
In [11]: import matplotlib.pyplot as plt
import seaborn as sns

# Sentiment distribution
sns.countplot(x='sentiment_category', data=data)
plt.title('Sentiment Distribution')
plt.show()

# Sentiment distribution by news category
plt.figure(figsize=(10, 6))
sns.countplot(x='Category', hue='sentiment_category', data=data)
plt.title('Sentiment Distribution by News Category')
plt.xticks(rotation=45)
plt.show()

```





## 5.4 News Title Language and Style Analysis

```
In [12]: from textstat import flesch_reading_ease

# Tokenization and lowercasing for frequency analysis
# Checking for non-string values before applying string methods
data['tokens'] = data['Title'].apply(lambda x: nltk.word_tokenize(x.lower()) if

# 1. Word Frequency Analysis
def word_frequency(tokens):
    freq_dist = FreqDist(tokens)
    return freq_dist.most_common(10)

data['word_freq'] = data['tokens'].apply(word_frequency)

# 2. N-gram Analysis (Bi-grams and Tri-grams)
def ngram_analysis(tokens, n=2):
    n_grams = ngrams(tokens, n)
    return FreqDist(n_grams).most_common(10)

# Bi-gram analysis
data['bigrams'] = data['tokens'].apply(lambda x: ngram_analysis(x, 2))

# Tri-gram analysis
data['trigrams'] = data['tokens'].apply(lambda x: ngram_analysis(x, 3))

# 3. Readability Analysis (Flesch Reading Ease)
def readability_score(title):
    if isinstance(title, str):
        return flesch_reading_ease(title)
    else:
        return None # or a default readability score
```

```
data['readability'] = data['Title'].apply(readability_score)

# Displaying the results
print(data[['word_freq', 'bigrams', 'trigrams', 'readability']].head())
```

```
word_freq \
DateTime
2023-08-30 06:54:49 [((, 3), (), 3), (,, 2), (chainlink, 1), (link...
2023-08-30 07:15:59 [(funds, 1), (punished, 1), (for, 1), (owning,...
2023-08-30 07:31:31 [(crude, 1), (oil, 1), (prices, 1), (stalled, ...
2023-08-30 10:38:40 [(grayscale, 1), ('s, 1), (bitcoin, 1), (win, ...
2023-08-30 11:00:00 [(i, 2), ('m, 2), (a, 1), (home, 1), (shopping...
```

```
bigrams \
DateTime
2023-08-30 06:54:49 [((chainlink, (), 1), (((, link), 1), ((link, ...
2023-08-30 07:15:59 [((funds, punished), 1), ((punished, for), 1),...
2023-08-30 07:31:31 [((crude, oil), 1), ((oil, prices), 1), ((pric...
2023-08-30 10:38:40 [((grayscale, 's), 1), (('s, bitcoin), 1), ((b...
2023-08-30 11:00:00 [((i, 'm), 2), (('m, a), 1), ((a, home), 1), (...]
```

```
trigrams \
DateTime
2023-08-30 06:54:49 [((chainlink, (, link), 1), (((, link, )), 1),...
2023-08-30 07:15:59 [((funds, punished, for), 1), ((punished, for,...
2023-08-30 07:31:31 [((crude, oil, prices), 1), ((oil, prices, sta...
2023-08-30 10:38:40 [((grayscale, 's, bitcoin), 1), (('s, bitcoin,...
2023-08-30 11:00:00 [((i, 'm, a), 1), (('m, a, home), 1), ((a, hom...
```

```
readability
DateTime
2023-08-30 06:54:49      76.22
2023-08-30 07:15:59      76.22
2023-08-30 07:31:31     113.10
2023-08-30 10:38:40      87.72
2023-08-30 11:00:00      73.17
```

Combine the traditional TF-IDF features with the new linguistic features (word frequency, n-grams, readability) to create an enhanced feature set. Use the combined feature set to train classification model.

```
In [13]: import numpy as np

from sklearn.preprocessing import MinMaxScaler
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import classification_report, accuracy_score
from scipy.sparse import hstack

# Convert word frequency lists into strings
data['word_freq_str'] = data['word_freq'].apply(lambda freqs: ' '.join(['_'.join(
data['bigrams_str'] = data['bigrams'].apply(lambda bigrams: ' '.join(['_'.join(
data['trigrams_str'] = data['trigrams'].apply(lambda trigrams: ' '.join(['_'.jo

# Vectorize the Word Frequencies, Bi-grams, and Tri-grams
vectorizer = CountVectorizer()
```

```

word_freq_vectorized = vectorizer.fit_transform(data['word_freq_str'])
bigrams_vectorized = vectorizer.fit_transform(data['bigrams_str'])
trigrams_vectorized = vectorizer.fit_transform(data['trigrams_str'])

# Normalize Readability Scores
scaler = MinMaxScaler()
data['readability_normalized'] = scaler.fit_transform(data[['readability']])

# Reshape 'readability_normalized' to match the number of rows
readability_normalized_2d = data['readability_normalized'].values.reshape(-1, 1)

from sklearn.impute import SimpleImputer

# Check for NaN values in each feature matrix
print("X_features NaN:", np.isnan(X_features.data).any())
print("word_freq_vectorized NaN:", np.isnan(word_freq_vectorized.data).any())
print("bigrams_vectorized NaN:", np.isnan(bigrams_vectorized.data).any())
print("trigrams_vectorized NaN:", np.isnan(trigrams_vectorized.data).any())
print("readability_normalized NaN:", np.isnan(readability_normalized_2d).any())

# Impute NaN values if necessary
# Impute NaN values in readability_normalized_2d with the mean
imputer = SimpleImputer(strategy='mean')
readability_normalized_imputed = imputer.fit_transform(readability_normalized_2d)

# Combine all features (use the imputed readability if NaN values were found)
combined_features = hstack([X_features, word_freq_vectorized, bigrams_vectorized,
                             readability_normalized_imputed])

# Train the Model
X_train, X_test, y_train, y_test = train_test_split(combined_features, data['Category'],
                                                    test_size=0.2, random_state=42)
model = MultinomialNB()
model.fit(X_train, y_train)

# Evaluate the model
y_pred = model.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Classification Report:\n", classification_report(y_test, y_pred))

```

X\_features NaN: False

word\_freq\_vectorized NaN: False

bigrams\_vectorized NaN: False

trigrams\_vectorized NaN: False

readability\_normalized NaN: True

Accuracy: 0.7279977140482128

Classification Report:

|               | precision | recall | f1-score | support |
|---------------|-----------|--------|----------|---------|
| Business      | 0.85      | 0.67   | 0.75     | 5197    |
| Entertainment | 0.90      | 0.75   | 0.82     | 5217    |
| Headlines     | 0.73      | 0.50   | 0.59     | 5468    |
| Health        | 0.85      | 0.77   | 0.81     | 3452    |
| Science       | 0.87      | 0.75   | 0.81     | 3297    |
| Sports        | 0.50      | 0.95   | 0.65     | 5825    |
| Technology    | 0.89      | 0.77   | 0.82     | 5020    |
| Worldwide     | 0.69      | 0.67   | 0.68     | 5020    |
| accuracy      |           |        | 0.73     | 38496   |
| macro avg     | 0.78      | 0.73   | 0.74     | 38496   |
| weighted avg  | 0.77      | 0.73   | 0.73     | 38496   |

## 6. Baseline Performance

In my project, I've established a baseline performance using a Naïve Bayes classifier, which is a commonly used method for text classification tasks due to its simplicity and effectiveness, especially in handling large datasets. The Naïve Bayes classifier serves as a robust starting point to compare more complex models against. This choice is supported by its frequent use in the literature for initial benchmarks in text classification problems.

For evaluating the baseline model, I employed standard metrics such as accuracy and the classification report, which includes precision, recall, and F1-score for each category. This approach allows for a comprehensive understanding of the model's performance across different news categories.

## 7. Classification Approach

For my classification task, I've chosen to utilize both the Multinomial Naïve Bayes and RandomForest classifiers. The selection of these classifiers is driven by their diverse strengths in handling text data. Multinomial Naïve Bayes is well-suited for large datasets and provides a good baseline due to its assumption of independence among features. On the other hand, RandomForest offers a more sophisticated approach by constructing multiple decision trees and merging their outputs. This method is beneficial for capturing more complex patterns in the data, which might be missed by the simpler Naïve Bayes model.

In terms of features, I've used TF-IDF vectorization for transforming the text data into a format suitable for machine learning models. TF-IDF helps in reflecting the importance of words in relation to the dataset as a whole, which is crucial for classification tasks like news categorization.

My contribution begins after the standard library implementations of these models and vectorization techniques. I've focused on integrating and fine-tuning these components, preprocessing the data, and thoroughly evaluating the models' performance to ensure the best possible results for the classification task.

# III Conclusions

## 9. Evaluation

In the evaluation of my classifier, I focused on using metrics such as accuracy, precision, recall, and F1-score. These metrics provide a comprehensive view of the model's performance, especially when dealing with imbalanced datasets, which is common in news categorization tasks.

- **Accuracy:** The accuracy of the classifier was around 73.48%, indicating a reasonable performance. However, this metric alone doesn't fully capture the model's effectiveness, especially in the context of class imbalances.



- **Precision, Recall, and F1-Score:** The detailed classification report highlighted the performance across different news categories. For example, the model showed high precision in categories like "Entertainment" and "Technology," indicating a low rate of false positives. However, categories like "Sports" had higher recall but lower precision, suggesting a tendency to over-classify articles in this category.

When compared to prior benchmarks, such as those reported in related literature, the performance of my approach is competitive, particularly given the complexity of classifying concise news titles. These benchmarks often report similar accuracy levels for basic classifiers, suggesting that my approach is on par with standard practices in the field.

## 10. Summary and Conclusions

Reflecting on the entirety of the project, the work stands out in its originality and ambition, particularly in tackling the less explored area of news title classification using advanced machine learning techniques. The novelty lies in the application of both Naïve Bayes and RandomForest classifiers, coupled with thorough feature engineering and preprocessing strategies tailored for the brevity and nuance of news headlines.

The project's contribution extends beyond academic interests; it offers practical applications in improving news categorization systems, which can enhance information retrieval and user experience in digital news platforms. Furthermore, the insights gained from the sentiment and trend analysis aspects of the project are valuable for understanding media patterns and public interests.

The methodologies and techniques employed are transferable and can be adapted to other text classification tasks, such as social media content analysis or document categorization in various digital libraries. The project's approach, being based on popular Python libraries and standard machine learning practices, can be replicated in different programming environments. This universality ensures that the methods can be leveraged by others in the field, potentially with different toolsets or for varied applications, thus contributing to the broader domain of natural language processing and machine learning.