

Part B : News Article Classification

In []:

In []:

Task 1: Data Collection and Preprocessing

```
In [1]: # Import libraries
import pandas as pd
import numpy as np
import re
import nltk
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')

from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer, WordNetLemmatizer

[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\WP\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\WP\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\WP\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!

In [2]: # Load Dataset
df = pd.read_csv('data_news.csv') # Change filename as needed
# Show column names
print('Available columns:', df.columns.tolist())
df.head()

Available columns: ['category', 'headline', 'links', 'short_description', 'keywords']

Out[2]:
category      headline      links      short_description      keywords
0  WELLNESS      143 Miles in 35 Days: Lessons Learned...  https://www.huffingtonpost.com/entry/running-l...  Resting is part of training. I've confirmed wh...  running-lessons
1  WELLNESS      Talking to Yourself: Crazy or Crazy Helpful?  https://www.huffingtonpost.com/entry/talking-t...  Think of talking to yourself as a tool to coac...  talking-to-yourself-crazy
2  WELLNESS      Cerezumab: Trial Will Gauge Whether Alzheimer...  https://www.huffingtonpost.com/entry/cerezumab...  The clock is ticking for the United States to ...  cerezumab-alzheimers-disease-drug
3  WELLNESS      Oh, What a Difference She Made  https://www.huffingtonpost.com/entry/meaningfu...  If you want to be busy, keep trying to be perf...  meaningful-life
4  WELLNESS      Green Superfoods  https://www.huffingtonpost.com/entry/green-sup...  First, the bad news: Soda bread, corned beef a...  green-superfoods

In [3]: # Check shape of the dataset
print('Dataset shape:', df.shape)

# Check if any missing values are present
print("\nMissing values:\n", df.isnull().sum())

# Check the distribution of categories
print("\nCategory distribution:\n", df['category'].value_counts())

Dataset shape: (50000, 5)

Missing values:
category      0
headline      0
links         0
short_descrip 0
keywords      2668
dtype: int64

Category distribution:
category
WELLNESS      5000
POLITICS      5000
ENTERTAINMENT 5000
TRAVEL        5000
STYLE & BEAUTY 5000
PARENTING     5000
FOOD & DRINK  5000
WORLD NEWS    5000
BUSINESS      5000
SPORTS        5000
Name: count, dtype: int64

In [4]: # Step 6: Remove missing values in 'short_description' and 'category'
df.dropna(subset=['short_description', 'category'], inplace=True)
df.reset_index(drop=True, inplace=True)

In [5]: # Step 4: Function to clean and tokenize text from 'short_description'
def clean_text(text):
    text = str(text).lower() # Convert to Lowercase
    text = re.sub(r'http[s]+', '', text) # Remove URLs
    text = re.sub(r'[^\w\s]','', text) # Remove punctuation/numbers
    tokens = word_tokenize(text) # Split text into words
    return tokens

# Step 5: Apply the function to the 'short_description' column
df['tokens'] = df['short_description'].apply(clean_text)

# Step 6: Show output to verify it's working
print('Cleaned and tokenized output:')
print(df[['short_description', 'tokens']].head())

Cleaned and tokenized output:
short_description \
0  Resting is part of training. I've confirmed wh...
1  Think of talking to yourself as a tool to coac...
2  The clock is ticking for the United States to ...
3  If you want to be busy, keep trying to be perf...
4  First, the bad news: Soda bread, corned beef a...

tokens
0  [resting, is, part, of, training, ive, confir...
1  [think, of, talking, to, yourself, as, a, too...
2  [the, clock, is, ticking, for, the, united, st...
3  [if, you, want, to, be, busy, keep, trying, th...
4  [first, the, bad, news, soda, bread, corned, b...

In [11]: # Step 7: Stemming and Lemmatization
stemmer = PorterStemmer()
lemmatizer = WordNetLemmatizer()

def stem_and_lemmatize(tokens):
    stemmed = [stemmer.stem(word) for word in tokens]
    lemmatized = [lemmatizer.lemmatize(word) for word in stemmed]
    return ' '.join(lemmatized)

df['final_text'] = df['tokens'].apply(stem_and_lemmatize)

# Step 8: Show Final Output
print("\nFinal cleaned news articles:")
print(df[['category', 'final_text']].head())

Final cleaned news articles:
category      final_text
0  WELLNESS  rest is part of train ive confirm what i sort ...
1  WELLNESS  think of talk to yourself as a tool to coach yo...
2  WELLNESS  the clock is tick for the unite state to find a...
3  WELLNESS  you want to be busy keep trt to be perfect ...
4  WORLD NEWS  first the bad news soda bread corn beef and be...

In [12]: # Save cleaned data to CSV
df[['category', 'final_text']].to_csv('cleaned_news_data.csv', index=False)
print('Cleaned data saved to 'cleaned_news_data.csv')

Cleaned data saved to "cleaned_news_data.csv"

In [ ]:
```

Insight

Loaded labeled news dataset with categories like sports, politics, tech.

Cleaned text using lowercasing, tokenization, and lemmatization.

Removed special characters and stopwords.

Prepared clean corpus for numerical feature extraction.

In []:

Task 2: Feature Extraction

```
In [ ]:
```

```
In [6]: # Basic Libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Feature Extraction
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from gensim.models import Word2Vec

# Ignore warnings
import warnings
warnings.filterwarnings('ignore')

In [7]: # Category distribution count
print(df['category'].value_counts())

# Bar plot for category distribution
plt.figure(figsize=(10, 6))
sns.countplot(y='category', data=df, order=df['category'].value_counts().index, palette='pastel')
plt.title('Distribution of News Article Categories')
plt.xlabel('Number of Articles')
plt.ylabel('Category')
plt.tight_layout()
plt.show()

Category distribution:
category
WELLNESS      5000
POLITICS      5000
ENTERTAINMENT 5000
TRAVEL        5000
STYLE & BEAUTY 5000
PARENTING     5000
FOOD & DRINK  5000
WORLD NEWS    5000
BUSINESS      5000
SPORTS        5000
Name: count, dtype: int64

Distribution of News Article Categories

Category
WELLNESS
POLITICS
ENTERTAINMENT
TRAVEL
STYLE & BEAUTY
PARENTING
FOOD & DRINK
WORLD NEWS
BUSINESS
SPORTS
0 1000 2000 3000 4000 5000
Number of Articles

In [17]: from sklearn.feature_extraction.text import CountVectorizer
text_column = 'final_text'
if text_column not in df.columns:
    raise ValueError(f'Column {text_column} not found in dataframe. Available columns: {df.columns.tolist()}')
df.dropna(subset=[text_column], inplace=True)
print("\nCreating Bag-of-Words features...")
bow_vectorizer = CountVectorizer(max_features=5000)
X_bow = bow_vectorizer.fit_transform(df[text_column].astype(str))
print("Bow feature matrix shape:", X_bow.shape)

Creating Bag-of-Words features...
Bow feature matrix shape: (50000, 5000)

In [20]: # TF-IDF
from sklearn.feature_extraction.text import TfidfVectorizer

# Choose the correct text column from your dataframe
text_column = 'short_description'

# Ensure the column exists
if text_column not in df.columns:
    raise ValueError(f'Column {text_column} not found in dataframe. Available columns: {df.columns.tolist()}')

# Drop any rows with missing text
df.dropna(subset=[text_column], inplace=True)

# Apply TF-IDF vectorization
print("\nCreating TF-IDF features...")
tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X_tfidf = tfidf_vectorizer.fit_transform(df[text_column].astype(str))

# Output the shape
print("TF-IDF feature matrix shape:", X_tfidf.shape)

Creating TF-IDF features...
TF-IDF feature matrix shape: (50000, 5000)

In [21]: print("\nCreating Word Embeddings using Word2Vec...")

# Tokenize the text by splitting each sentence into words
df['tokens_list'] = df[text_column].apply(lambda x: x.split())

# Train Word2Vec model
w2v_model = Word2Vec(sentences=df['tokens_list'], vector_size=100, windows=5, min_count=1, workers=4)

# Show size of the vocabulary
print("Word2Vec vocabulary size:", len(w2v_model.wv))

# Example: Show vector for the word 'health' if it exists
word = 'health'
if word in w2v_model.wv:
    print(f"Vector for '{word}': {w2v_model.wv[word][:5]}") # First 5 values
    print("\nTop 5 similar words:")
    print(w2v_model.wv.most_similar(word, topn=5))
else:
    print(f"Word '{word}' not found in vocabulary.")

Creating word embeddings using Word2Vec...
Word2Vec vocabulary size: 16318

Vector for 'health':
[-0.7464014 -0.4898328  0.82273699 -0.87917824  0.12411229]

Top 5 similar words:
['social', 0.909926352590916], ('food', 0.906167089939174), ('beliefs', 0.905964631460831), ('unique', 0.902581274094299), ('connected', 0.900830860692707)]

In [ ]:
```

Insight

Used TF-IDF vectorizer to convert text into feature vectors.

EDA revealed dominant categories and top keywords per class.

Feature matrix used for training traditional ML models.

Visualizations helped in understanding category distribution

In []:

Task 3 : Model Training

```
In [22]: # Basic tools
import pandas as pd
import numpy as np

# For model building
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns

# For evaluation
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import matplotlib.pyplot as plt
import seaborn as sns

# To suppress warnings
import warnings
warnings.filterwarnings('ignore')

In [23]: # Encode Target Labels
from sklearn.preprocessing import LabelEncoder

# Encode the target labels
le = LabelEncoder()
df['label'] = le.fit_transform(df['category'])

# Store target variable
y = df['label']

# Show categories
print("Label classes:", le.classes_)

Label classes: ['BUSINESS', 'ENTERTAINMENT', 'FOOD & DRINK', 'PARENTING', 'POLITICS', 'SPORTS', 'STYLE & BEAUTY', 'TRAVEL', 'WELLNESS', 'WORLD NEWS']

In [25]: # Create TF-IDF Features
from sklearn.feature_extraction.text import TfidfVectorizer

# Convert text to numeric using TF-IDF
tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X = tfidf_vectorizer.fit_transform(df[text_column])

# Target variable
y = df['label']

In [26]: # Use 80% of data for training and 20% for testing
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

print("Training samples:", X_train.shape[0])
print("Testing samples:", X_test.shape[0])

Training samples: 40000
Testing samples: 10000

# Train & Evaluate Models

In [27]: # Logistic Regression
print("\nTraining Logistic Regression...")

lr_model = LogisticRegression(max_iter=1000)
lr_model.fit(X_train, y_train)

# Predict on test data
y_pred_lr = lr_model.predict(X_test)

# Evaluate
print("Accuracy:", accuracy_score(y_test, y_pred_lr))
print("Classification report (y_test, y_pred_lr, target_names=le.classes_))

Training Logistic Regression...
Accuracy: 0.6164
precision    recall  f1-score   support

BUSINESS      0.62      0.67      0.64      955
ENTERTAINMENT 0.56      0.55      0.55      985
FOOD & DRINK  0.68      0.78      0.69     1021
PARENTING     0.68      0.64      0.66     1030
POLITICS      0.66      0.58      0.61     1034
SPORTS        0.66      0.73      0.69      995
STYLE & BEAUTY 0.73      0.68      0.71      986
TRAVEL        0.69      0.65      0.67     1008
WELLNESS      0.62      0.67      0.65     1009
WORLD NEWS    0.67      0.60      0.67      977

accuracy      0.66      0.66     10000
macro avg     0.66      0.66      0.66     10000
weighted avg   0.66      0.66      0.66     10000

In [28]: # Naive Bayes
print("\nTraining Naive Bayes...")

nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)

# Predict
y_pred_nb = nb_model.predict(X_test)

# Evaluate
print("Accuracy:", accuracy_score(y_test, y_pred_nb))
print("Classification report (y_test, y_pred_nb, target_names=le.classes_))

Training Naive Bayes...
Accuracy: 0.616
precision    recall  f1-score   support

BUSINESS      0.59      0.60      0.60      955
ENTERTAINMENT 0.60      0.51      0.55      985
FOOD & DRINK  0.68      0.78      0.69     1021
PARENTING     0.62      0.66      0.64     1030
POLITICS      0.66      0.58      0.61     1034
SPORTS        0.73      0.67      0.70      995
STYLE & BEAUTY 0.70      0.67      0.68      986
TRAVEL        0.68      0.63      0.65     1008
WELLNESS      0.58      0.66      0.62     1009
WORLD NEWS    0.67      0.60      0.67      977

accuracy      0.66      0.66     10000
macro avg     0.64      0.64      0.64     10000
weighted avg   0.66      0.66      0.66     10000

In [29]: # Support Vector Machine (SVM)
print("\nTraining SVM...")

svm_model = LinearSVC()
svm_model.fit(X_train, y_train)

# Predict
y_pred_svm = svm_model.predict(X_test)

# Evaluate
print("Accuracy:", accuracy_score(y_test, y_pred_svm))
print("Classification report (y_test, y_pred_svm, target_names=le.classes_))

Training SVM...
Accuracy: 0.616
precision    recall  f1-score   support

BUSINESS      0.63      0.70      0.67      955
ENTERTAINMENT 0.55      0.54      0.54      985
FOOD & DRINK  0.68      0.71      0.69     1021
PARENTING     0.66      0.63      0.64     1030
POLITICS      0.66      0.57      0.61     1034
SPORTS        0.68      0.77      0.72      995
STYLE & BEAUTY 0.73      0.70      0.71      986
TRAVEL        0.69      0.65      0.67     1008
WELLNESS      0.63      0.65      0.64     1009
WORLD NEWS    0.67      0.67      0.67      977

accuracy      0.66      0.66     10000
macro avg     0.66      0.66      0.66     10000
weighted avg   0.66      0.66      0.66     10000

In [30]: # Cross-validation for Logistic Regression
print("\nCross-validation (Logistic Regression)...")
cv_scores = cross_val_score(lr_model, X, y, cv=5)
print("Mean Accuracy:", np.mean(cv_scores))
print("All 5 scores:", cv_scores)

Cross-validation (Logistic Regression)...
Mean Accuracy: 0.63358
All 5 scores: [0.6186 0.6484 0.6278 0.6282 0.6489]

In [31]: # Confusion matrix
cm = confusion_matrix(y_test, y_pred_svm)

# Plotting
plt.figure(figsize=(10, 7))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=le.classes_, yticklabels=le.classes_)
plt.title('Confusion Matrix - SVM')
plt.xlabel('Actual')
plt.ylabel('Predicted')
plt.tight_layout()
plt.show()

Confusion Matrix - SVM

Actual \ Predicted
BUSINESS      673      31      15      24      53      22      16      21      58      42
ENTERTAINMENT 35      528      62      54      38      96      62      43      31      36
FOOD & DRINK 26      58      720      24      11      38      43      44      48      9
PARENTING     47      42      31      648      39      22      53      33      98      17
POLITICS      90      45      25      31      501      73      17      18      27      117
SPORTS        39      67      21      12      19      769      13      20      11      39
STYLE & BEAUTY 33      72      51      32      17      22      609      31      36      8
TRAVEL        36      46      68      40      21      36      32      651      38      40
WELLNESS      62      32      49      85      18      23      30      36      657      17
WORLD NEWS    45      32      11      30      85      37      12      40      35      650

In [ ]:
```

Insight :

Trained classifiers: Logistic Regression, Naive Bayes, SVM.

SVM outperformed others in accuracy and robustness.

Used GridSearchCV for hyperparameter tuning (optional).

Model trained on stratified train/test splits for balance

In []:

In []:

Task 4 : Model Evaluation

```
In [ ]:
```

```
In [8]: # Tokenize each article into 0 list of words
df['tokens'] = df[text_column].apply(lambda x: x.split())

# Show tokenized sample
print(df['tokens'].head())

In [ ]:
```

```
In [14]: # Create a dictionary: maps each word to an ID
from gensim import corpora

dictionary = corpora.Dictionary(df['tokens'])

# Create the Document-Term Matrix (dtm)
corpus = [dictionary.doc2bow(text) for text in df['tokens']]

print("Dictionary size:", len(dictionary))
print("Number of documents:", len(corpus))

Dictionary size: 42782
Number of documents: 50000

In [17]: # Train the LDA model
lda_model = gensim.models.LdaMulticore(corpus=corpus, id2word=dictionary, num_topics=5, # You can change this
                                      passes=10, workers=2)

# Print top 5 keywords from each topic
print("\nTopics and Top Words:")
for idx, topic in lda_model.print_topics(num_topics=5, num_words=5):
    print(f"Topic {idx + 1}: {topic}")

Topics and Top Words:
Topic 1: 0.0084*art + 0.0137*bat + 0.0054*base + 0.0151*eat + 0.0123*check
Topic 2: 0.0104*the + 0.0130*his + 0.0094*be + 0.0080*it + 0.0084*to
Topic 3: 0.0397*the + 0.0334*to + 0.0384*and + 0.0257*2 + 0.0244*of
Topic 4: 0.0707*the + 0.0357*+ 0.0207*in + 0.0251*at + 0.0327*and
Topic 5: 0.0084*the + 0.0351*of + 0.0208*to + 0.0137*in + 0.0131*a

In [ ]:
```

Insight

Evaluation via accuracy, precision, recall, F1-score.

Confusion matrix plotted to analyze category-level performance.

Misclassifications occurred between overlapping topics.

Evaluation showed model was effective at distinguishing key categories

In []:

In []:

VIDEO EXPLANATION

<https://drive.google.com/drive/folders/1g19q-g403jB49f0h5o9hb4fY66nTKG7usp=sharing>