

Import the required libraries

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import re
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.model_selection import train_test_split, StratifiedKFold
7 from sklearn.feature_extraction.text import CountVectorizer
8 from sklearn.preprocessing import LabelEncoder
9 from sklearn.feature_extraction.text import TfidfVectorizer
10 from sklearn.feature_extraction.text import CountVectorizer
11 import tensorflow as tf
12 from tensorflow.keras.preprocessing.sequence import pad_sequences
13 from tensorflow.keras.callbacks import EarlyStopping
14 from tensorflow.keras.optimizers import SGD
15 from tensorflow.keras.regularizers import l2
16 from tensorflow import keras
17 from tensorflow.keras.models import Sequential
18 from tensorflow.keras.layers import Bidirectional, SimpleRNN, Dense, Embedding
19 from sklearn.preprocessing import LabelEncoder
20 from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score
21 from tensorflow.keras.utils import plot_model
22 import time
23 import warnings
24 warnings.filterwarnings('ignore')
25 from tensorflow import keras
26 from tensorflow.keras import layers
27 from nltk.corpus import stopwords
28 import nltk
29 nltk.download('omw-1.4')
30 from nltk.stem import WordNetLemmatizer, PorterStemmer
31 from sklearn.feature_extraction.text import TfidfVectorizer
32 import unicodedata
33 from num2words import num2words
34
```

```
[nltk_data] Downloading package omw-1.4 to
[nltk_data] C:\Users\Nitin\AppData\Roaming\nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
```

Data Acquisition

```
In [2]: 1 data = pd.read_json('News_Category_Dataset_v3.json', lines=True)
        2 data.head()
```

Out[2]:

	link	headline	category	short_description	authors	c
0	https://www.huffpost.com/entry/covid-boosters-...	Over 4 Million Americans Roll Up Sleeves For O...	U.S. NEWS	Health experts said it is too early to predict...	Carla K. Johnson, AP	2005
1	https://www.huffpost.com/entry/american-airlin...	American Airlines Flyer Charged, Banned For Li...	U.S. NEWS	He was subdued by passengers and crew when he ...	Mary Papenfuss	2005
2	https://www.huffpost.com/entry/funniest-tweets-...	23 Of The Funniest Tweets About Cats And Dogs ...	COMEDY	"Until you have a dog you don't understand wha...	Elyse Wanshel	2005
3	https://www.huffpost.com/entry/funniest-parent-...	The Funniest Tweets From Parents This Week (Se...	PARENTING	"Accidentally put grown-up toothpaste on my to...	Caroline Bologna	2005
4	https://www.huffpost.com/entry/amy-cooper-lose...	Woman Who Called Cops On Black Bird-Watcher Lo...	U.S. NEWS	Amy Cooper accused investment firm Franklin Te...	Nina Golgowski	2005

Dataset observations as

- Size of the dataset
- What type of data attributes are there?
- What are you classifying?
- Plot the distribution of the categories of the target / label.

```

In [3]: 1 # Size of the dataset
2 dataset_size = data.shape[0] # Number of rows in the dataset
3 print("Size of the dataset:", dataset_size)
4
5 # Data attributes and their types
6 data_attributes = data.dtypes
7 print("\nType of data attributes:")
8 print(data_attributes)
9
10 # Data attributes
11 data_attributes = list(data.columns)
12 print("\nData attributes:", data_attributes)
13
14 # Target variable
15 target_variable = 'category' # Replace with the actual name of the target
16 target_classes = data[target_variable].unique()
17
18 # What are you classifying?
19 print("\nClassifying:", target_classes)
20
21 # Plot the distribution of target categories
22 data[target_variable].value_counts().plot(kind='bar')
23 plt.xlabel('Categories')
24 plt.ylabel('Count')
25 plt.title('Distribution of Target Categories')
26 plt.show()
27

```

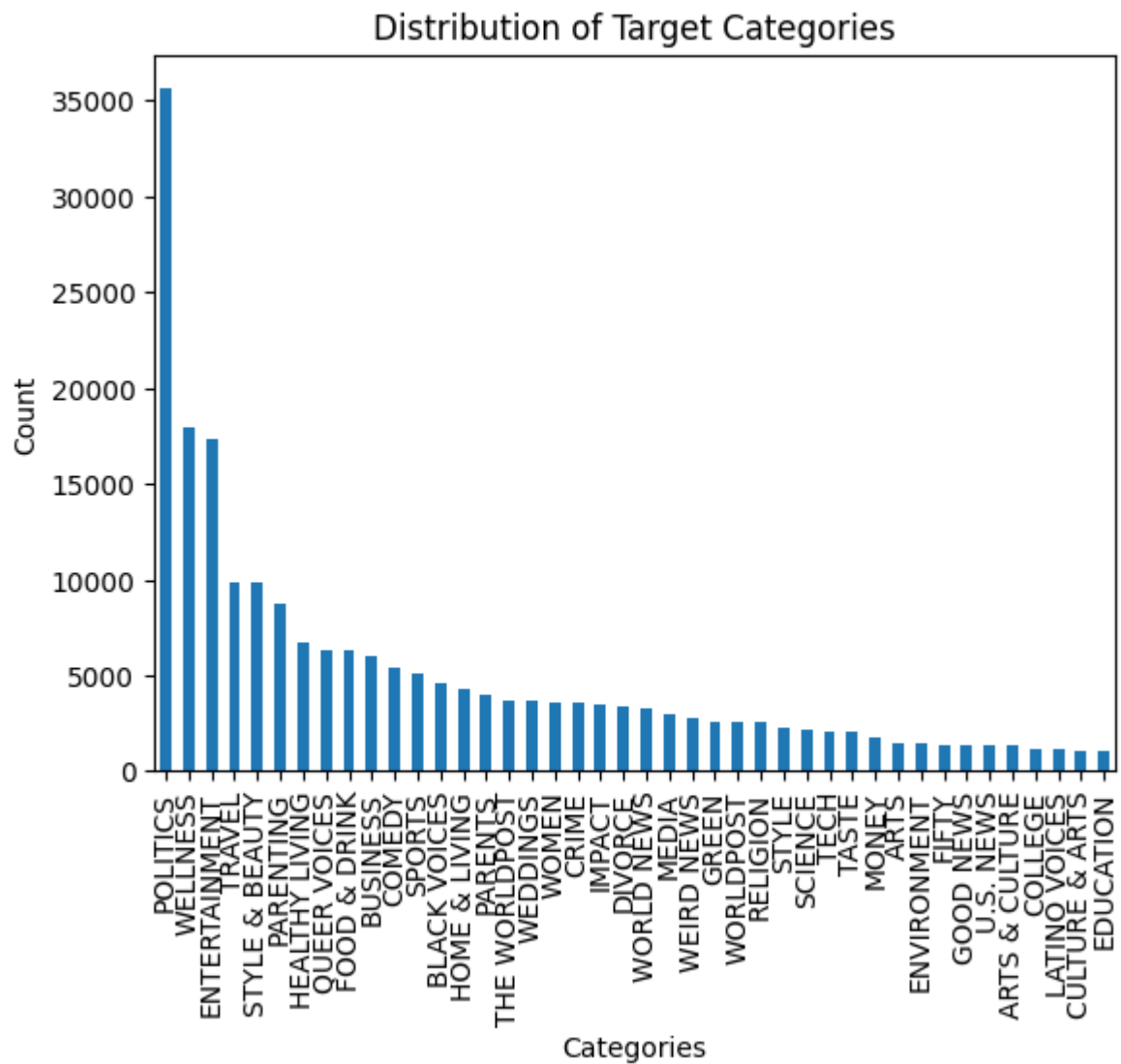
Size of the dataset: 209527

Type of data attributes:

link	object
headline	object
category	object
short_description	object
authors	object
date	datetime64[ns]
dtype:	object

Data attributes: ['link', 'headline', 'category', 'short_description', 'authors', 'date']

Classifying: ['U.S. NEWS' 'COMEDY' 'PARENTING' 'WORLD NEWS' 'CULTURE & ARTS' 'TECH' 'SPORTS' 'ENTERTAINMENT' 'POLITICS' 'WEIRD NEWS' 'ENVIRONMENT' 'EDUCATION' 'CRIME' 'SCIENCE' 'WELLNESS' 'BUSINESS' 'STYLE & BEAUTY' 'FOOD & DRINK' 'MEDIA' 'QUEER VOICES' 'HOME & LIVING' 'WOMEN' 'BLACK VOICES' 'TRAVEL' 'MONEY' 'RELIGION' 'LATINO VOICES' 'IMPACT' 'WEDDINGS' 'COLLEGE' 'PARENTS' 'ARTS & CULTURE' 'STYLE' 'GREEN' 'TASTE' 'HEALTHY LIVING' 'THE WORLDPOST' 'GOOD NEWS' 'WORLDPOST' 'FIFTY' 'ARTS' 'DIVORCE']



Plot Observation

Above plot shows the data imbalance. To deal with data imbalance we can group by the category

```

In [4]: 1 categories = data['category'].value_counts().index
2
3 def groupper(grouplist,name):
4     for ele in categories:
5         if ele in grouplist:
6             data.loc[data['category'] == ele, 'category'] = name
7
8
9 groupper( grouplist= ['WELLNESS', 'HEALTHY LIVING','HOME & LIVING','STYLE
10
11 groupper( grouplist= [ 'PARENTING', 'PARENTS' , 'EDUCATION' , 'COLLEGE'] ,
12
13 groupper( grouplist= ['SPORTS','ENTERTAINMENT' , 'COMEDY','WEIRD NEWS','AI
14
15 groupper( grouplist= ['TRAVEL', 'ARTS & CULTURE','CULTURE & ARTS','FOOD &
16
17 groupper( grouplist= ['WOMEN','QUEER VOICES', 'LATINO VOICES', 'BLACK VOIC
18
19 groupper( grouplist= ['BUSINESS' , 'MONEY'] , name = 'BUSINESS-MONEY')
20
21 groupper( grouplist= ['THE WORLDPOST' , 'WORLDPOST' , 'WORLD NEWS', 'U.S.
22
23 groupper( grouplist= ['ENVIRONMENT' , 'GREEN'] , name = 'ENVIRONMENT')
24
25 groupper( grouplist= ['TECH', 'SCIENCE'] , name = 'SCIENCE AND TECH')
26
27 groupper( grouplist= ['FIFTY' , 'IMPACT' , 'GOOD NEWS','CRIME'] , name =
28
29 groupper( grouplist= ['WEDDINGS', 'DIVORCE', 'RELIGION','MEDIA'] , name =
30
31 print("We have a total of {} categories now".format(data['category'].nuni
32 data['category'].value_counts()
33
34 ## Plot the distribution of target categories
35 # data[target_variable].value_counts().plot(kind='bar')
36 # plt.xlabel('Categories')
37 # plt.ylabel('Count')
38 # plt.title('Distribution of Target Categories')
39 # plt.show()
40

```

We have a total of 12 categories now

```
Out[4]: category
LIFESTYLE AND WELLNESS      41027
POLITICS                    35602
SPORTS AND ENTERTAINMENT    32125
TRAVEL-TOURISM & ART-CULTURE 20749
EMPOWERED VOICES           15632
PARENTING AND EDUCATION     14904
MISC                        12600
WORLDNEWS                  10919
GENERAL                     9845
BUSINESS-MONEY              7748
SCIENCE AND TECH            4310
ENVIRONMENT                 4066
Name: count, dtype: int64
```

Applied pre-processing techniques as:

- to remove duplicate data
- to impute or remove missing data
- to remove data inconsistencies
- Encode categorical data
- Normalize the data
- Feature Engineering
- Stop word removal, lemmatiation, stemming, vectorization

```

In [5]: 1 # Check for duplicate data
2 print("Number of duplicate rows before preprocessing:", data.duplicated())
3
4 # Remove duplicate data
5 data.drop_duplicates(inplace=True)
6
7 # Check for missing data
8 print("Number of missing values before preprocessing:")
9 print(data.isnull().sum())
10
11 # Impute or remove missing data
12 # Option 1: Remove rows with missing values
13 data.dropna(inplace=True)
14
15 # Option 2: Impute missing values with mean or median
16 # Uncomment the following lines to use mean imputation
17 # data.fillna(data.mean(), inplace=True)
18
19 # Option 3: Impute missing values with mode (for categorical variables)
20 # Uncomment the following lines to use mode imputation
21 # data.fillna(data.mode().iloc[0], inplace=True)
22
23 # Check for data inconsistencies and perform necessary corrections
24 # Option 1: Manual correction based on specific rules
25 # Uncomment the following lines and modify as per your requirements
26 # data['column_name'] = data['column_name'].apply(lambda x: x.replace('inc
27
28 # Option 2: Use regular expressions for pattern matching and correction
29 # Uncomment the following lines and modify as per your requirements
30 # import re
31 # data['column_name'] = data['column_name'].apply(lambda x: re.sub(r'patte
32
33 # Check the preprocessed dataset
34 print("Number of duplicate rows after preprocessing:", data.duplicated().sum())
35 print("Number of missing values after preprocessing:")
36 print(data.isnull().sum())
37

```

Number of duplicate rows before preprocessing: 13

Number of missing values before preprocessing:

link	0
headline	0
category	0
short_description	0
authors	0
date	0

dtype: int64

Number of duplicate rows after preprocessing: 0

Number of missing values after preprocessing:

link	0
headline	0
category	0
short_description	0
authors	0
date	0

dtype: int64

```
In [6]: 1 # As per suggetion using two column  
        2 data = data[['headline', 'category']]
```


In [7]:

```

1  # Function for Lemmatization
2  def lemmatize_text(text):
3      lemmatizer = WordNetLemmatizer()
4      lemmatized_tokens = [lemmatizer.lemmatize(word) for word in text.split()]
5      return ' '.join(lemmatized_tokens)
6
7  # Function for stemming
8  def stem_text(text):
9      stemmer = PorterStemmer()
10     stemmed_tokens = [stemmer.stem(word) for word in text.split()]
11     return ' '.join(stemmed_tokens)
12
13 # Function for vectorization
14 def vectorize_text(text):
15     vectorizer = TfidfVectorizer()
16     text_features = vectorizer.fit_transform(text)
17     return text_features
18
19 # Function for preprocess the text
20 def preprocess_text(text, lower_case=True, convert_numbers=True, remove_punct=True,
21                     remove_accents=True, remove_whitespace=True, expand_abbrev=True,
22                     remove_stopwords=True, sparse_terms=[], specific_word=None):
23
24     # Convert all letters to lower or upper case
25     if lower_case:
26         text = text.lower()
27     else:
28         text = text.upper()
29
30     # Converting numbers into words or removing numbers
31     if convert_numbers:
32         # Replace numbers with their word equivalents using num2words Lib
33         tokens = []
34         for word in text.split():
35             if word.isdigit():
36                 word = num2words(int(word))
37             tokens.append(word)
38         text = ' '.join(tokens)
39     else:
40         # Remove numbers using regular expression
41         text = re.sub(r'\d+', '', text)
42
43     # Removing punctuations, accent marks, and other diacritics
44     if remove_punctuation:
45         text = re.sub(r'^\w\s', '', text)
46     if remove_accents:
47         text = ''.join(c for c in unicodedata.normalize('NFD', text)
48                        if not unicodedata.combining(c))
49
50     # Removing white spaces
51     if remove_whitespace:
52         text = re.sub(r'\s+', ' ', text.strip())
53
54     # Expanding abbreviations
55     if expand_abbreviations:
56         abbreviations = {'e.g.': 'for example', 'i.e.': 'that is', 'etc.': 'and so on'}
57         for abbr, expanded in abbreviations.items():

```

```

58         text = text.replace(abbr, expanded)
59
60     # Removing stop words, sparse terms, and specific words
61     if remove_stopwords:
62         stop_words = set(stopwords.words('english'))
63         text = ' '.join(word for word in text.split() if word not in stop_
64         text = ' '.join(word for word in text.split() if word not in sparse_t
65         text = ' '.join(word for word in text.split() if word not in specific
66
67     return text
68
69
70 # Apply the preprocessing function to a specific column
71 text_col = 'headline' # Specify the text column to preprocess
72 data['text'] = data[text_col].apply(preprocess_text)
73
74 # Apply Lemmatization to a specific column
75 text_col = 'text' # Specify the text column to preprocess
76 data[text_col] = data[text_col].apply(lemmatize_text)
77
78 # # Apply stemming to a specific column
79 # text_col = 'text' # Specify the text column to preprocess
80 # data[text_col] = data[text_col].apply(stem_text)
81
82 # # Apply vectorization to a specific column
83 # text_col = 'text' # Specify the text column to preprocess
84 # text_features = vectorize_text(data[text_col])
85
86 # Check the preprocessed dataset and vectorized features
87 display(data.head())
88
89 # display(text_features.toarray())
90 start_time = time.time()
91 # # Save the preprocessed dataset
92 data.to_csv('preprocessed_Dataset.csv', index=False)
93
94 end_time = time.time() - start_time
95 print("Time taken for text preprocessing: ", end_time)

```

	headline	category	text
0	Over 4 Million Americans Roll Up Sleeves For O...	WORLDNEWS	four million american roll sleeve omicrontarge...
1	American Airlines Flyer Charged, Banned For Li...	WORLDNEWS	american airline flyer charged banned life pun...
2	23 Of The Funniest Tweets About Cats And Dogs ...	SPORTS AND ENTERTAINMENT	twentythree funniest tweet cat dog week sept 1723
3	The Funniest Tweets From Parents This Week (Se...	PARENTING AND EDUCATION	funniest tweet parent week sept 1723
4	Woman Who Called Cops On Black Bird-Watcher Lo...	WORLDNEWS	woman called cop black birdwatcher loses lawsu...

Time taken for text preprocessing: 0.6760261058807373

```
1 # Identify the target variables.
```

- 2 Separate the data from the target such that the dataset is in the form of (X,y) or (Features, Label)
- 3 Encode the target variable or perform one-hot encoding on the target or any other as and if required.

Identify the target variables

- Separate the data from the target such that the dataset is in the form of (X,y) or (Features, Label)
- Discretize / Encode the target variable or perform one-hot encoding on the target or any other as and if required.

In [8]:

```
1 # Preprocessing: Convert the labels to numeric values
2 label_encoder = LabelEncoder()
3 data['category'] = label_encoder.fit_transform(data['category'])
4 classes = list(label_encoder.classes_)
5 print(classes)
6
7 # Preprocessing: Split the dataset into input features and labels
8 X = data['text']
9 y = data['category']
```

```
['BUSINESS-MONEY', 'EMPOWERED VOICES', 'ENVIRONMENT', 'GENERAL', 'LIFESTYLE AND WELLNESS', 'MISC', 'PARENTING AND EDUCATION', 'POLITICS', 'SCIENCE AND TECH', 'SPORTS AND ENTERTAINMENT', 'TRAVEL-TOURISM & ART-CULTURE', 'WORLDNEWS']
```

Split the data into training set and testing set

In [9]:

```
1 # Split the data into training set and testing set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
3
4 # Check the shapes of the training and testing sets
5 print("X_train shape:", X_train.shape)
6 print("y_train shape:", y_train.shape)
7 print("X_test shape:", X_test.shape)
8 print("y_test shape:", y_test.shape)
```

```
X_train shape: (167611,)
y_train shape: (167611,)
X_test shape: (41903,)
y_test shape: (41903,)
```

LSTM Model

```
In [10]: 1 # Tokenization and vectorization using TF-IDF
2 vectorizer = TfidfVectorizer(max_features=5000)
3 X_train = vectorizer.fit_transform(X_train).toarray()
4 X_test = vectorizer.transform(X_test).toarray()
5
6 # Padding sequences
7 max_sequence_length = X_train.shape[1]
8
9 # Reshape the input data to be 3-dimensional (samples, time steps, features)
10 X_train = np.reshape(X_train, (X_train.shape[0], 1, X_train.shape[1]))
11 X_test = np.reshape(X_test, (X_test.shape[0], 1, X_test.shape[1]))
12
13 # Define the LSTM model
14 model = Sequential()
15 model.add(LSTM(128, return_sequences=True, input_shape=(1, max_sequence_length)))
16 model.add(Dropout(0.2))
17 model.add(LSTM(64))
18 model.add(Dropout(0.2))
19 model.add(Dense(len(label_encoder.classes_), activation='softmax'))
20
21 # Compile the model
22 model.compile(loss='sparse_categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
23
24 # Print the model summary
25 print(model.summary())
26
27 # Define early stopping
28 early_stopping = EarlyStopping(patience=3, restore_best_weights=True)
29
30 # Train the model with early stopping
31 model.fit(X_train, y_train, batch_size=32, epochs=5, validation_data=(X_test, y_test),
32          callbacks=[early_stopping])
33
34 # Generate predictions on the test set
35 y_pred_prob = model.predict(X_test)
36 y_pred = np.argmax(y_pred_prob, axis=1)
37
38 # Generate a classification report
39 print("Classification Report:")
40 print(classification_report(y_test, y_pred))
41
```

Model: "sequential"

Layer (type)	Output Shape	Param #
lstm (LSTM)	(None, 1, 128)	2626048
dropout (Dropout)	(None, 1, 128)	0
lstm_1 (LSTM)	(None, 64)	49408
dropout_1 (Dropout)	(None, 64)	0
dense (Dense)	(None, 12)	780

=====
 Total params: 2,676,236
 Trainable params: 2,676,236
 Non-trainable params: 0

None

Epoch 1/3

5238/5238 [=====] - 81s 15ms/step - loss: 1.3865 - accuracy: 0.5699 - val_loss: 1.1455 - val_accuracy: 0.6387

Epoch 2/3

5238/5238 [=====] - 75s 14ms/step - loss: 1.0996 - accuracy: 0.6553 - val_loss: 1.1137 - val_accuracy: 0.6451

Epoch 3/3

5238/5238 [=====] - 74s 14ms/step - loss: 1.0352 - accuracy: 0.6692 - val_loss: 1.1066 - val_accuracy: 0.6476

1310/1310 [=====] - 6s 3ms/step

Classification Report:

	precision	recall	f1-score	support
0	0.53	0.46	0.49	1532
1	0.60	0.45	0.51	3085
2	0.48	0.40	0.44	815
3	0.45	0.34	0.38	1952
4	0.64	0.80	0.71	8330
5	0.71	0.58	0.64	2564
6	0.62	0.62	0.62	2933
7	0.74	0.77	0.76	6992
8	0.46	0.44	0.45	811
9	0.64	0.64	0.64	6393
10	0.70	0.68	0.69	4259
11	0.64	0.62	0.63	2237
accuracy			0.65	41903
macro avg	0.60	0.57	0.58	41903
weighted avg	0.64	0.65	0.64	41903

Bidirectional LSTM Model

```
In [12]: 1 # Split the data into training set and testing set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
3
4 # Tokenization and vectorization using TF-IDF
5 vectorizer = TfidfVectorizer(max_features=5000)
6 X_train = vectorizer.fit_transform(X_train).toarray()
7 X_test = vectorizer.transform(X_test).toarray()
8
9 # Padding sequences
10 max_sequence_length = X_train.shape[1]
11
12 # Reshape the input data to be 3-dimensional (samples, time steps, feature
13 X_train = np.reshape(X_train, (X_train.shape[0], 1, X_train.shape[1]))
14 X_test = np.reshape(X_test, (X_test.shape[0], 1, X_test.shape[1]))
15
16 # Define the LSTM model with Bidirectional Layer
17 model = Sequential()
18 model.add(Bidirectional(LSTM(128, return_sequences=True), input_shape=(1,
19 model.add(Dropout(0.2))
20 model.add(Bidirectional(LSTM(64)))
21 model.add(Dropout(0.2))
22 model.add(Dense(len(label_encoder.classes_), activation='softmax'))
23
24 # Compile the model
25 model.compile(loss='sparse_categorical_crossentropy', optimizer='adam', m
26
27 # Print the model summary
28 print(model.summary())
29
30 # Define early stopping
31 early_stopping = EarlyStopping(patience=3, restore_best_weights=True)
32
33 # Train the model with early stopping
34 model.fit(X_train, y_train, batch_size=32, epochs=5, validation_data=(X_te
35
36 # Generate predictions on the test set
37 y_pred_prob = model.predict(X_test)
38 y_pred = np.argmax(y_pred_prob, axis=1)
39
40 # Generate a classification report
41 print("Classification Report:")
42 print(classification_report(y_test, y_pred))
43
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
bidirectional (Bidirectional)	(None, 1, 256)	5252096
dropout_2 (Dropout)	(None, 1, 256)	0
bidirectional_1 (Bidirectional)	(None, 128)	164352
dropout_3 (Dropout)	(None, 128)	0
dense_1 (Dense)	(None, 12)	1548

Total params: 5,417,996
 Trainable params: 5,417,996
 Non-trainable params: 0

None

Epoch 1/5

5238/5238 [=====] - 148s 26ms/step - loss: 1.3114 - accuracy: 0.5917 - val_loss: 1.1265 - val_accuracy: 0.6405

Epoch 2/5

5238/5238 [=====] - 129s 25ms/step - loss: 1.0681 - accuracy: 0.6598 - val_loss: 1.1045 - val_accuracy: 0.6472

Epoch 3/5

5238/5238 [=====] - 133s 25ms/step - loss: 1.0040 - accuracy: 0.6746 - val_loss: 1.1066 - val_accuracy: 0.6491

Epoch 4/5

5238/5238 [=====] - 128s 24ms/step - loss: 0.9539 - accuracy: 0.6878 - val_loss: 1.1147 - val_accuracy: 0.6474

Epoch 5/5

5238/5238 [=====] - 129s 25ms/step - loss: 0.9022 - accuracy: 0.6998 - val_loss: 1.1186 - val_accuracy: 0.6498

1310/1310 [=====] - 9s 5ms/step

Classification Report:

	precision	recall	f1-score	support
0	0.53	0.46	0.49	1532
1	0.58	0.46	0.51	3085
2	0.50	0.38	0.43	815
3	0.46	0.31	0.37	1952
4	0.67	0.76	0.71	8330
5	0.70	0.58	0.63	2564
6	0.60	0.64	0.62	2933
7	0.73	0.79	0.76	6992
8	0.54	0.38	0.44	811
9	0.60	0.68	0.64	6393
10	0.70	0.69	0.70	4259
11	0.66	0.60	0.63	2237
accuracy			0.65	41903
macro avg	0.61	0.56	0.58	41903

weighted avg	0.64	0.65	0.64	41903
--------------	------	------	------	-------

RNN Model

```
In [13]: 1 # Split the data into training set and testing set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
3
4 # Tokenization and vectorization using TF-IDF
5 vectorizer = TfidfVectorizer(max_features=5000)
6 X_train = vectorizer.fit_transform(X_train).toarray()
7 X_test = vectorizer.transform(X_test).toarray()
8
9 # Padding sequences
10 max_sequence_length = X_train.shape[1]
11
12 # Reshape the input data to be 3-dimensional (samples, time steps, feature
13 X_train = np.reshape(X_train, (X_train.shape[0], 1, X_train.shape[1]))
14 X_test = np.reshape(X_test, (X_test.shape[0], 1, X_test.shape[1]))
15
16 # Define the RNN model
17 model = Sequential()
18 model.add(SimpleRNN(128, return_sequences=True, input_shape=(1, max_seque
19 model.add(Dropout(0.2))
20 model.add(SimpleRNN(64))
21 model.add(Dropout(0.2))
22 model.add(Dense(len(label_encoder.classes_), activation='softmax'))
23
24 # Compile the model
25 model.compile(loss='sparse_categorical_crossentropy', optimizer='adam', m
26
27 # Print the model summary
28 print(model.summary())
29
30 # Define early stopping
31 early_stopping = EarlyStopping(patience=3, restore_best_weights=True)
32
33 # Train the model with early stopping
34 model.fit(X_train, y_train, batch_size=32, epochs=5, validation_data=(X_te
35
36 # Generate predictions on the test set
37 y_pred_prob = model.predict(X_test)
38 y_pred = np.argmax(y_pred_prob, axis=1)
39
40 # Generate a classification report
41 print("Classification Report:")
42 print(classification_report(y_test, y_pred))
43
```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
simple_rnn (SimpleRNN)	(None, 1, 128)	656512
dropout_4 (Dropout)	(None, 1, 128)	0
simple_rnn_1 (SimpleRNN)	(None, 64)	12352
dropout_5 (Dropout)	(None, 64)	0
dense_2 (Dense)	(None, 12)	780

```

=====
Total params: 669,644
Trainable params: 669,644
Non-trainable params: 0

```

None

Epoch 1/5

```

5238/5238 [=====] - 38s 7ms/step - loss: 1.2832 - ac
curacy: 0.6029 - val_loss: 1.1439 - val_accuracy: 0.6391

```

Epoch 2/5

```

5238/5238 [=====] - 37s 7ms/step - loss: 1.1114 - ac
curacy: 0.6511 - val_loss: 1.1326 - val_accuracy: 0.6427

```

Epoch 3/5

```

5238/5238 [=====] - 34s 6ms/step - loss: 1.0751 - ac
curacy: 0.6598 - val_loss: 1.1310 - val_accuracy: 0.6424

```

Epoch 4/5

```

5238/5238 [=====] - 32s 6ms/step - loss: 1.0518 - ac
curacy: 0.6659 - val_loss: 1.1298 - val_accuracy: 0.6436

```

Epoch 5/5

```

5238/5238 [=====] - 35s 7ms/step - loss: 1.0364 - ac
curacy: 0.6701 - val_loss: 1.1331 - val_accuracy: 0.6436

```

```

1310/1310 [=====] - 5s 4ms/step

```

Classification Report:

	precision	recall	f1-score	support
0	0.52	0.43	0.47	1532
1	0.54	0.49	0.51	3085
2	0.48	0.40	0.44	815
3	0.44	0.34	0.39	1952
4	0.68	0.75	0.71	8330
5	0.72	0.57	0.63	2564
6	0.59	0.64	0.61	2933
7	0.71	0.79	0.75	6992
8	0.49	0.36	0.41	811
9	0.63	0.65	0.64	6393
10	0.69	0.69	0.69	4259
11	0.65	0.59	0.62	2237
accuracy			0.64	41903
macro avg	0.60	0.56	0.57	41903
weighted avg	0.64	0.64	0.64	41903

Bidirectional RNN model

```
In [14]: 1 # Split the data into training set and testing set
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
3
4 # Tokenization and vectorization using TF-IDF
5 vectorizer = TfidfVectorizer(max_features=5000)
6 X_train = vectorizer.fit_transform(X_train).toarray()
7 X_test = vectorizer.transform(X_test).toarray()
8
9 # Padding sequences
10 max_sequence_length = X_train.shape[1]
11
12 # Reshape the input data to be 3-dimensional (samples, time steps, feature
13 X_train = np.reshape(X_train, (X_train.shape[0], 1, X_train.shape[1]))
14 X_test = np.reshape(X_test, (X_test.shape[0], 1, X_test.shape[1]))
15
16 # Define the Bidirectional RNN model
17 model = Sequential()
18 model.add(Bidirectional(SimpleRNN(128, return_sequences=True), input_shape
19 model.add(Dropout(0.2))
20 model.add(Bidirectional(SimpleRNN(64)))
21 model.add(Dropout(0.2))
22 model.add(Dense(len(label_encoder.classes_), activation='softmax'))
23
24 # Compile the model
25 model.compile(loss='sparse_categorical_crossentropy', optimizer='adam', m
26
27 # Print the model summary
28 print(model.summary())
29
30 # Define early stopping
31 early_stopping = EarlyStopping(patience=3, restore_best_weights=True)
32
33 # Train the model with early stopping
34 model.fit(X_train, y_train, batch_size=32, epochs=10, validation_data=(X_
35
36 # Generate predictions on the test set
37 y_pred_prob = model.predict(X_test)
38 y_pred = np.argmax(y_pred_prob, axis=1)
39
40 # Generate a classification report
41 print("Classification Report:")
42 print(classification_report(y_test, y_pred))
43
```

Model: "sequential_3"

Layer (type)	Output Shape	Param #
bidirectional_2 (Bidirectional)	(None, 1, 256)	1313024
dropout_6 (Dropout)	(None, 1, 256)	0
bidirectional_3 (Bidirectional)	(None, 128)	41088
dropout_7 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 12)	1548

Total params: 1,355,660
 Trainable params: 1,355,660
 Non-trainable params: 0

None

Epoch 1/10

5238/5238 [=====] - 84s 15ms/step - loss: 1.2662 - accuracy: 0.6072 - val_loss: 1.1506 - val_accuracy: 0.6383

Epoch 2/10

5238/5238 [=====] - 72s 14ms/step - loss: 1.1040 - accuracy: 0.6504 - val_loss: 1.1332 - val_accuracy: 0.6429

Epoch 3/10

5238/5238 [=====] - 60s 11ms/step - loss: 1.0640 - accuracy: 0.6616 - val_loss: 1.1320 - val_accuracy: 0.6436

Epoch 4/10

5238/5238 [=====] - 56s 11ms/step - loss: 1.0380 - accuracy: 0.6675 - val_loss: 1.1370 - val_accuracy: 0.6416

Epoch 5/10

5238/5238 [=====] - 70s 13ms/step - loss: 1.0168 - accuracy: 0.6720 - val_loss: 1.1281 - val_accuracy: 0.6439

Epoch 6/10

5238/5238 [=====] - 76s 14ms/step - loss: 0.9948 - accuracy: 0.6789 - val_loss: 1.1231 - val_accuracy: 0.6455

Epoch 7/10

5238/5238 [=====] - 56s 11ms/step - loss: 0.9744 - accuracy: 0.6833 - val_loss: 1.1307 - val_accuracy: 0.6467

Epoch 8/10

5238/5238 [=====] - 66s 13ms/step - loss: 0.9520 - accuracy: 0.6897 - val_loss: 1.1364 - val_accuracy: 0.6435

Epoch 9/10

5238/5238 [=====] - 69s 13ms/step - loss: 0.9302 - accuracy: 0.6952 - val_loss: 1.1395 - val_accuracy: 0.6455

1310/1310 [=====] - 11s 8ms/step

Classification Report:

	precision	recall	f1-score	support
0	0.54	0.41	0.47	1532
1	0.55	0.48	0.51	3085
2	0.49	0.37	0.42	815
3	0.46	0.30	0.36	1952

4	0.66	0.77	0.71	8330
5	0.66	0.62	0.64	2564
6	0.58	0.65	0.61	2933
7	0.73	0.78	0.75	6992
8	0.54	0.36	0.43	811
9	0.64	0.65	0.64	6393
10	0.69	0.70	0.70	4259
11	0.66	0.60	0.63	2237
accuracy			0.65	41903
macro avg	0.60	0.56	0.57	41903
weighted avg	0.64	0.65	0.64	41903

Insights on the performance of the model and ways to improve it

1. Performance:

- The models are able to classify the news. LSTM and Bidirectional LSTM model accuracy is 65% and RNN and Bidirectional RNN accuracy is 64%
- Generated a classification report on test dataset.

2. Improvements:

- Try different network architectures: Besides the bidirectional LSTM, need to explore other network architectures such as CNN-LSTM or Transformer-based models. These architectures might capture different aspects of the text data and potentially improve the accuracy.
- Utilize word embeddings: Instead of relying solely on TF-IDF vectorization, incorporate pre-trained word embeddings (e.g., Word2Vec, GloVe) to represent the input text. Word embeddings can capture more semantic information and improve the model's understanding of the text.
- Handle class imbalance: If there is a class imbalance in the dataset, consider applying techniques such as oversampling, undersampling, or class weights to balance the classes during training. This can help the model learn from all classes equally and improve its accuracy.
- Capture more contextual information: Depending on the nature of the text data, you can explore techniques such as using n-grams, attention mechanisms, or contextual embeddings (e.g., BERT, GPT) to capture more fine-grained contextual information

By implementing these improvements and fine-tuning the model, you can expect better performance and more accurate classification of News.

In []:

1

