# BSc Computer Science

# CM3060 Natural Language Processing

# Comparative Analysis of Statistical and Embedding-Based Models for News Categorization

# 1. Introduction

This coursework endeavors to develop a text classifier applied to the domain of fake news detection, utilizing both statistical and embedding-based language models. The objective is to comprehensively compare the efficacy of traditional statistical methods against modern deep learning approaches in addressing this critical problem.

## 1.1 Domain-Specific Area

The exponential growth of online content has made it increasingly challenging to manually categorize and sift through vast amounts of news articles available on the internet. Automated systems are crucial to efficiently classify news and provide users with relevant information quickly and accurately. This research addresses the need for automated news classification using machine learning models and Natural Language Processing (NLP), focusing on predicting news categories based on news titles.

Automated news classification systems play a pivotal role in modern information retrieval by categorizing news articles into predefined topics or themes. This classification not only aids in organizing information but also facilitates personalized content recommendations tailored to individual user interests. By automating this process, users can access pertinent news promptly, enhancing their overall browsing experience and information consumption efficiency.

Despite significant advancements in NLP, automatically categorizing news headlines remains a complex task due to varying news lengths and nuanced semantics. Traditional approaches like word count or frequency-based methods, such as count vectorizers and TF-IDF, have limitations in capturing semantic meanings embedded in news titles. These methods often rely on lexical features and may not effectively discern the underlying context of news topics.

In recent years, machine learning models have emerged as powerful tools for text classification tasks. Statistical models like Naïve Bayes, logistic regression, and ensemble methods such as decision trees and random forests have been widely employed for their ability to handle large datasets and achieve high accuracy in classification tasks.

Furthermore, deep learning techniques leveraging convolutional and recurrent neural networks have shown promise in learning complex features and patterns from textual data.

To address the shortcomings of traditional methods, embedding-based models like Word2Vec and FastText have gained prominence. Word2Vec constructs vector representations of words that capture semantic similarities, enhancing classification features derived from lexical features like bag-of-words and TF-IDF. FastText, on the other hand, considers word internal structures and is particularly effective for morphologically rich languages, improving the robustness of word representations.

## 1.2 Research Objectives

The primary objective of this project is to explore and compare the effectiveness of both statistical and embedding-based models in the task of news categorization, specifically focusing on news titles.The goal is to assess how well traditional statistical models, such as Naïve Bayes and logistic regression using count vectorizers and TF-IDF, perform in comparison to modern deep learning models like LSTM, and embedding-based models like Word2Vec. This evaluation will provide insights into which types of models are most effective for the task of news categorization based on title information.

## 1.3 Literature Review

This section provides a comprehensive review of significant studies relevant to feature engineering and machine learning classifiers in the context of Natural Language Processing (NLP) and automated news classification.

### 1.3.1 Feature Engineering

Feature engineering plays a crucial role in developing effective NLP systems, as it involves transforming raw data into meaningful features that facilitate machine learning algorithms. In the realm of text classification, preserving the context of textual data while reducing dimensionality is essential for accurate classification.

**i) Frequency-based Word Embedding**

Traditional approaches such as Bag of Words (BoW), term frequency, and TF-IDF are foundational in representing words within the vector space model, making them computationally efficient. However, these methods often neglect the semantic relationships between words and fail to capture contextual nuances.

In a study inspired by previous work, a supervised classification method for news articles achieved 95% accuracy by analyzing titles and constructing inputs based on single words and variable-sized n-grams. Despite its high accuracy, this method's feature extraction process lacks the ability to preserve intricate semantic relationships within documents. This research explores methods akin to those discussed in for feature extraction and emphasizes leveraging semantic features derived from advanced embedding techniques.

Comparative studies on word embedding techniques across different domains, including biomedical and social media analyses, highlight the efficiency and accuracy of prediction-based embeddings like Word2Vec and FastText. These embeddings excel in capturing semantic meanings and are adept at contextualizing words within a document.

**ii) Prediction-based Word Embedding**

Mikolov et al. introduced Word2Vec, employing two architectures—continuous bag-of-words (CBOW) and skip-gram—to learn and represent words in vector spaces. Their experiments demonstrated superior accuracy in semantic and syntactic tasks with increased dimensional size and vocabulary, albeit limited to local context within documents.

FastText, an extension of Word2Vec proposed by Stein et al., addresses the limitation of Word2Vec by considering each word as composed of character n-grams. This approach allows FastText to create embeddings that incorporate morphological information, making it particularly effective for morphologically rich languages. Experiments on datasets like IMDB have shown FastText achieving significant precision improvements, emphasizing its capability to capture complex word structures and enhance classification accuracy.

## 1.3.2 Classification

Classification of news articles based on extracted features is a fundamental task in automated news categorization systems. Various machine learning models, particularly statistical classifiers, have been employed to establish robust relationships between news titles and predefined categories.

Al-Tahrawi et al. utilized logistic regression (LR) to classify news categories, achieving high precision rates, notably 96.5% for one category and over 90% for three out of five categories. Similarly, Naïve Bayes (NBC) has been applied effectively in news article classification, demonstrating average recall values of 92.87% and precision values of 91.16% in classifying diverse news categories.

Research by Kabir et al. compared classifiers such as Stochastic Gradient Descent (SGD) against NBC and LR for Bangla news titles, highlighting SGD's superior performance. Hybrid models combining different classifiers and data sampling techniques have also been explored, with Pambudi et al. achieving high accuracy in multi-class classification of Indonesian news using Pseudo Nearest Neighbor (PNNR) methods.

Furthermore, integrating advanced embeddings such as Word2Vec and FastText into traditional classification models has shown promising results in enhancing semantic understanding and improving classification accuracy.

The reviewed literature underscores the pivotal role of feature engineering and machine learning classifiers in advancing automated news classification systems. While traditional methods like frequency-based embeddings and statistical classifiers provide robust

performance, prediction-based embeddings such as Word2Vec and FastText offer significant improvements in capturing semantic nuances within news articles.

## 1.4 Dataset Description

The dataset used for this research is sourced from the "News Aggregator" dataset available in the UCI Machine Learning Repository, collected from March 10 to August 10 of 2014. This dataset comprises references to news web pages gathered from an online aggregator, organized into clusters that group pages discussing the same news story. The dataset also includes references to web pages that link to one of the news pages within the collection.

Content and Format: The dataset is provided in tab-delimited CSV format.

**uci-news-aggregator.csv**

- **Columns:**
  - `ID` : Numeric ID of the news page
  - `TITLE` : Title of the news article
  - `URL` : URL of the news page
  - `PUBLISHER` : Name of the publisher
  - `CATEGORY` : News category ( `b` for business, `t` for science and technology, `e` for entertainment, `m` for health)
  - `STORY` : Alphanumeric ID of the cluster containing news about the same story
  - `HOSTNAME` : Hostname of the URL
  - `TIMESTAMP` : Approximate publication time of the news article in milliseconds since January 1, 1970, GMT

I have downloaded the dataset from kaggle.

Link to Dataset: https://www.kaggle.com/datasets/uciml/news-aggregator-dataset

The dataset is utilized for tasks such as classification and clustering of news articles based on their titles and associated metadata. It provides a comprehensive collection of news articles across multiple categories, enabling researchers to explore various NLP and machine learning techniques for automated news categorization.

Source and Citation: Gasparetti, F. (2016). News Aggregator. UCI Machine Learning Repository. Retrieved from https://doi.org/10.24432/C5F61C.

## 1.5 Evaluation Methodology

In this section, we outline the evaluation methodology used to assess the performance of the text classification models, employing both statistical and embedding-based approaches. We discuss the metrics utilized to measure the effectiveness of each methodology and how these metrics facilitate a comparative analysis.

When evaluating classification models in this research, several metrics are employed to gauge their performance across different aspects:

1. **Accuracy** : Measures the ratio of correctly predicted instances to the total number of instances evaluated. Accuracy provides an overall assessment of how well the model predicts across all classes, suitable for balanced datasets.

2. **Precision** : Precision measures the proportion of true positive predictions among all positive predictions made by the model. It is crucial in scenarios where the cost of false positives is high, emphasizing the model's ability to correctly identify relevant instances.

3. **Recall** : Recall calculates the proportion of true positive predictions among all instances that actually belong to the positive class. This metric is particularly important in contexts where missing positive instances (false negatives) are costly, ensuring comprehensive coverage of relevant cases.

4. **F1-Score** : The F1-Score is the harmonic mean of precision and recall, providing a balanced measure that considers both false positives and false negatives. It is useful when there is an uneven class distribution or when both precision and recall need to be equally weighted in evaluation.

5. **Cohen's Kappa** : Cohen's kappa measures the agreement between two annotators on a classification problem, correcting for the agreement occurring by chance. This metric is valuable when evaluating models in scenarios where class distributions are not uniform or when assessing inter-rater reliability.

6. **Balanced Accuracy** : Balanced accuracy calculates the average recall obtained in each class, useful in handling imbalanced datasets. It provides a fair representation of model performance across all classes, even when class sizes vary significantly.

7. **Matthews Correlation Coefficient (MCC)** : MCC considers true and false positives and negatives, offering a balanced measure of classification quality. It is particularly robust in scenarios with imbalanced class distributions and is widely used for binary and multiclass classifications.

By applying a comprehensive set of evaluation metrics, this research aims to provide a thorough comparison between statistical and embedding-based approaches in the context of automated news classification.

# 2. Data Preparation

In this phase, we begin by loading and exploring the 'News Aggregator' dataset.

## 2.1 Installing Libraries and Importing Dependencies

```python
# Import necessary libraries
import numpy as np
import pandas as pd
import re
import string

# Data visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Machine Learning models and evaluation metrics
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import SGDClassifier, LogisticRegression
from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearc
from sklearn_evaluation.plot import confusion_matrix
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score
from sklearn.model_selection import cross_val_score
from sklearn.multiclass import OneVsRestClassifier
from sklearn.feature_selection import chi2
from sklearn.metrics import cohen_kappa_score, hamming_loss, log_loss, zero_one_los

# Natural Language Processing tools
import nltk
from nltk.tokenize import sent_tokenize, word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

# Deep Learning with Keras and TensorFlow
import keras.backend as K
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, LSTM, Dense
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.preprocessing.text import Tokenizer

# Word Embeddings with Gensim
import gensim
from gensim.models import Word2Vec, KeyedVectors

# Additional utility libraries
from tabulate import tabulate
import collections
import gzip
import shutil

# Download necessary NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')

# Set maximum display width for pandas DataFrame
pd.set_option('display.max_colwidth', 100)
```

```
# Suppress warnings
import warnings
warnings.filterwarnings("ignore")

# Set up plotting styles for matplotlib and seaborn
sns.set(style="whitegrid")
plt.style.use('ggplot')

# Clear TensorFlow session to avoid clutter from previous runs
K.clear_session()

# Ensure specific version of scipy is installed (if needed)
# !pip install scipy==1.12

# Additional setup for specific tasks or models can be added here
```

```
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\lavan\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\lavan\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\lavan\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
WARNING:tensorflow:From C:\Users\lavan\anaconda3\Lib\site-packages\keras\src\backend
\common\global_state.py:82: The name tf.reset_default_graph is deprecated. Please us
e tf.compat.v1.reset_default_graph instead.
```

## 2.2 Loading and Initial Data Exploration

Load the dataset, display the head, and count the number of entries and category counts.

In [2]:
```
# Load the dataset
news = pd.read_csv('data/uci-news-aggregator.csv')
```

In [3]:
```
# Display the first few rows of the dataset
print("First few rows of the dataset:")
news.head()
```

```
First few rows of the dataset:
```

Out[3]:

| | ID | TITLE | URL | PUBLISHER |
|---|---|---|---|---|
| 0 | 1 | Fed official says weak data caused by weather, should not slow taper | http://www.latimes.com/business/money/la-fi-mo-federal-reserve-plosser-stimulus-economy-20140310... | Los Angeles Times |
| 1 | 2 | Fed's Charles Plosser sees high bar for change in pace of tapering | http://www.livemint.com/Politics/H2EvwJSK2VE6OF7iK1g3PP/Feds-Charles-Plosser-sees-high-bar-for-c... | Livemint |
| 2 | 3 | US open: Stocks fall after Fed official hints at accelerated tapering | http://www.ifamagazine.com/news/us-open-stocks-fall-after-fed-official-hints-at-accelerated-tape... | IFA Magazine |
| 3 | 4 | Fed risks falling 'behind the curve', Charles Plosser says | http://www.ifamagazine.com/news/fed-risks-falling-behind-the-curve-charles-plosser-says-294430 | IFA Magazine |
| 4 | 5 | Fed's Plosser: Nasty Weather Has Curbed Job Growth | http://www.moneynews.com/Economy/federal-reserve-charles-plosser-weather-job-growth/2014/03/10/i... | Moneynews |

In [4]:
```
# Count the number of non-null values in each column
print("Count of non-null values per column:")
news.count()
```

Count of non-null values per column:

```
Out[4]:  ID            422419
         TITLE         422419
         URL           422419
         PUBLISHER     422417
         CATEGORY      422419
         STORY         422419
         HOSTNAME      422419
         TIMESTAMP     422419
         dtype: int64
```

Each column's count is printed, showing the number of non-null entries for the columns.

```
In [5]:  # Count the occurrences of each category in the 'CATEGORY' column
         print("\nCounts of each category in the 'CATEGORY' column:")
         news.CATEGORY.value_counts()
```

```
         Counts of each category in the 'CATEGORY' column:
```

```
Out[5]:  CATEGORY
         e    152469
         b    115967
         t    108344
         m     45639
         Name: count, dtype: int64
```

The result is a Series where the index represents each unique category ('e', 'b', 't', 'm') and the values represent the count of each category.

```
In [6]:  # Reset seaborn style to default
         sns.set(style='whitegrid')

         # Get the counts of each category in the 'CATEGORY' column
         category_counts = news['CATEGORY'].value_counts()

         # Define a mapping for category labels to human-readable names
         cat_map = {
             'b': 'Business',
             't': 'Science',
             'e': 'Entertainment',
             'm': 'Health'
         }

         # Sort the categories based on the defined order for plotting
         # Create lists of categories and their corresponding counts in the desired order
         categories = [cat_map[cat] for cat in sorted(cat_map.keys())]
         counts = [category_counts[cat] for cat in sorted(cat_map.keys())]

         # Plotting
         plt.figure(figsize=(8, 6))
         sns.barplot(x=categories, y=counts, palette='viridis')
         plt.xlabel('Category')
         plt.ylabel('Count')
         plt.title('Category Counts')
         plt.xticks(rotation=45)
         plt.show()
```

Category Counts

# 3. Data Preprocessing

In this phase, we will preprocess the data for effective analysis and modeling, ensuring that the dataset is clean, normalized, and balanced for accurate classification of news articles based on their content.

## 3.1 Balancing and Preparing Categorical Labels for Classification

This code snippet demonstrates the process of balancing the 'News Aggregator' dataset by selecting an equal number of samples from each of the four categories (business, entertainment, science, health). After combining and shuffling the selected samples, numerical labels are assigned to each category for classification purposes.

In [7]:
```python
# Define the number of samples per category to balance the dataset
num_of_categories = 45639

# Shuffle the original dataset
shuffled = news.reindex(np.random.permutation(news.index))
```

```python
# Select equal numbers of samples from each category for balanced dataset
e = shuffled[shuffled['CATEGORY'] == 'e'][:num_of_categories]
b = shuffled[shuffled['CATEGORY'] == 'b'][:num_of_categories]
t = shuffled[shuffled['CATEGORY'] == 't'][:num_of_categories]
m = shuffled[shuffled['CATEGORY'] == 'm'][:num_of_categories]

# Concatenate the selected samples into a single DataFrame
concated = pd.concat([e, b, t, m], ignore_index=True)

# Shuffle the concatenated dataset
concated = concated.reindex(np.random.permutation(concated.index))

# Assign numerical labels based on categories for classification
concated['LABEL'] = 0
concated.loc[concated['CATEGORY'] == 'e', 'LABEL'] = 0
concated.loc[concated['CATEGORY'] == 'b', 'LABEL'] = 1
concated.loc[concated['CATEGORY'] == 't', 'LABEL'] = 2
concated.loc[concated['CATEGORY'] == 'm', 'LABEL'] = 3
```

In [8]:
```python
# Display the first 10 labels and their corresponding categorical representations
print(concated['LABEL'][:10])
labels = to_categorical(concated['LABEL'], num_classes=4)
print(labels[:10])
```

```
79647     1
140189    3
48375     1
47025     1
69394     1
177560    3
55413     1
25491     0
4909      0
110271    2
Name: LABEL, dtype: int64
[[0. 1. 0. 0.]
 [0. 0. 0. 1.]
 [0. 1. 0. 0.]
 [0. 1. 0. 0.]
 [0. 1. 0. 0.]
 [0. 0. 0. 1.]
 [0. 1. 0. 0.]
 [1. 0. 0. 0.]
 [1. 0. 0. 0.]
 [0. 0. 1. 0.]]
```

The output includes the first 10 labels and their categorical representations, along with the counts of entries in the balanced dataset and the distribution of categories post-balancing. This ensures that the dataset is ready

In [9]:
```python
# Drop the 'CATEGORY' column if it exists in the DataFrame
if 'CATEGORY' in concated.keys():
    concated.drop(['CATEGORY'], axis=1)
```

```
In [10]:  # Count the number of entries in each column of the balanced dataset
          concated.count()

Out[10]:  ID           182556
          TITLE        182556
          URL          182556
          PUBLISHER    182555
          CATEGORY     182556
          STORY        182556
          HOSTNAME     182556
          TIMESTAMP    182556
          LABEL        182556
          dtype: int64

In [11]:  # Display the count of each category after balancing
          concated['CATEGORY'].value_counts()

Out[11]:  CATEGORY
          b    45639
          m    45639
          e    45639
          t    45639
          Name: count, dtype: int64

In [12]:  # Get the counts of each category in the 'CATEGORY' column
          category_counts = concated['LABEL'].value_counts()

          # Define the label map
          label_map = {
              0: 'Entertainment',
              1: 'Business',
              2: 'Science',
              3: 'Health'
          }

          # Sort the categories based on the defined order
          labels = [label_map[label] for label in sorted(label_map.keys())]
          counts = [category_counts[label] for label in sorted(label_map.keys())]

          # Plotting
          plt.figure(figsize=(8, 6))
          sns.barplot(x=labels, y=counts, palette='viridis')
          plt.xlabel('Category')
          plt.ylabel('Count')
          plt.title('Category Counts After Balancing')
          plt.xticks(rotation=45)
          plt.show()
```

Category Counts After Balancing

## 3.2 Preprocessing for Word Embeddings

The approach in this research adopted is similar to the baseline model discussed in the literature but there are several additional changes implemented that provided bette results.

### 3.2.1 Data cleaning - Removing blanks and duplicates

The code removes duplicate rows from the DataFrame `concated`, considering all columns. It then checks for null values in each column and prints the number of null values per column.

In [13]:
```python
# Remove duplicates based on all columns
concated = concated.drop_duplicates()

# Check for null values in each column
concated.isnull().sum()
```

```
Out[13]: ID           0
         TITLE        0
         URL          0
         PUBLISHER    1
         CATEGORY     0
         STORY        0
         HOSTNAME     0
         TIMESTAMP    0
         LABEL        0
         dtype: int64
```

In [14]: `concated.info()`

```
<class 'pandas.core.frame.DataFrame'>
Index: 182556 entries, 79647 to 65397
Data columns (total 9 columns):
 #   Column     Non-Null Count   Dtype
---  ------     --------------   -----
 0   ID         182556 non-null  int64
 1   TITLE      182556 non-null  object
 2   URL        182556 non-null  object
 3   PUBLISHER  182555 non-null  object
 4   CATEGORY   182556 non-null  object
 5   STORY      182556 non-null  object
 6   HOSTNAME   182556 non-null  object
 7   TIMESTAMP  182556 non-null  int64
 8   LABEL      182556 non-null  int64
dtypes: int64(3), object(6)
memory usage: 13.9+ MB
```

Finally, `concated.info()` provides a summary of the DataFrame, showing non-null counts and data types for each column.

### 3.2.2 Concatenating columns - Title & Publisher

Combination of 'Title' column with 'Publisher' in a new column 'Text.

In [15]:
```python
# Concatenate 'TITLE' and 'PUBLISHER' into 'text'
concated['text'] = concated['TITLE'] + " " + concated['PUBLISHER']

# Confirm the changes
concated.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 182556 entries, 79647 to 65397
Data columns (total 10 columns):
 #   Column     Non-Null Count    Dtype
---  ------     --------------    -----
 0   ID         182556 non-null   int64
 1   TITLE      182556 non-null   object
 2   URL        182556 non-null   object
 3   PUBLISHER  182555 non-null   object
 4   CATEGORY   182556 non-null   object
 5   STORY      182556 non-null   object
 6   HOSTNAME   182556 non-null   object
 7   TIMESTAMP  182556 non-null   int64
 8   LABEL      182556 non-null   int64
 9   text       182555 non-null   object
dtypes: int64(3), object(7)
memory usage: 15.3+ MB
```

In [16]:
```python
# Verify concatenation
concated.text.iloc[50]  # Check a sample row
```

Out[16]:  "Preview 'True Blood' Season 7 Episode 3 'Fire in the Hole' Geeks of Doom"

### 3.2.3 Removing punctuation, stopwords, lowercasing, and tokenization

This data pre-processing step removes punctuation, stopwords, lowercasing, and generates tokens.

In [17]:
```python
symbols = [',', '.', '"', ':', ')', '(', '-', '!', '?', '|',
          ';', "'", '$', '&', '/', '[', ']', '>', '%', '=',
          '#', '*', '+', '\\', '•', '~', '@', '£', '·', '_',
          '{', '}', '©', '^', '®', '`', '<', '→', '°', '€',
          '™', '›', '♥', '←', '×', '§', '″', '′', 'Â', '█',
          '½', 'à', '…', '"', '★', '"', '–', '●', 'â', '►',
          '−', '¢', '²', '¬', '░', '¶', '↑', '±', '¿', '▾',
          '═', '¦', '║', '―', '¥', '▓', '—', '‹', '─', '▒', '：',
          '¼', '⊕', '▼', '▪', '†', '■', '’', '▀', '¨', '▄', '♫',
          '☆', 'é', '¯', '♦', '¤', '▲', 'è', '¸', '¾', 'Ã', '⋅',
          '‘', '∞', '∙', ')', '↓', '、', '│', '(', '»', '，', '♪',
          '╩', '╚', '³', '・', '╦', '╣', '╔', '╗', '▬', '❤', 'ï', 'Ø',
          '¹', '≤', '‡', '√', ]
```

In [18]:
```python
# Defines a function clean_symbol to remove specific symbols from text and applies
def clean_symbol(text):
    text = str(text)
    for symbol in symbols:
        text = text.replace(symbol, '')
    return text

# remove symbols and punctuations
concated['text'] = concated['text'].apply(lambda x: clean_symbol(x))
```

```python
In [19]:  #remove all non-word characters (punctuations) from the 'text' column
          concated['text'] = concated['text'].str.replace('[^\w\s]','')
```

```python
In [20]:  # Defines a function remove numbers and extra spaces from text and applies it to th
          def clean_text(s):
              s = re.sub("[0-9]+", "",s)
              s = re.sub(' +',' ', s)
              return s

          concated['text'] = [clean_text(s) for s in concated['text']]
```

```python
In [21]:  # Confirm the changes
          concated.text.iloc[50]
```

```
Out[21]:  'Preview True Blood Season Episode Fire in the Hole Geeks of Doom'
```

```python
In [22]:  # Dictionary of short form words and mispellings
          short_forms_dict = {"ain't": "is not", "aren't": "are not","can't": "cannot",
                              "'cause": "because", "could've": "could have", "couldn't": "cou
                              "didn't": "did not",  "doesn't": "does not", "don't": "do not",
                              "hasn't": "has not", "haven't": "have not", "he'd": "he would",
                              "he's": "he is", "how'd": "how did", "how'd'y": "how do you", "
                              "how's": "how is",  "I'd": "I would", "I'd've": "I would have",
                              "I'll've": "I will have","I'm": "I am", "I've": "I have", "i'd"
                              "i'd've": "i would have", "i'll": "i will",  "i'll've": "i will
                              "i've": "i have", "isn't": "is not", "it'd": "it would", "it'd'
                              "it'll": "it will", "it'll've": "it will have","it's": "it is",
                              "ma'am": "madam", "mayn't": "may not", "might've": "might have"
                              "mightn't've": "might not have", "must've": "must have", "mustn
                              "mustn't've": "must not have", "needn't": "need not", "needn't'
                              "o'clock": "of the clock", "oughtn't": "ought not", "oughtn't'v
                              "shan't": "shall not", "sha'n't": "shall not", "shan't've": "sh
                              "she'd": "she would", "she'd've": "she would have", "she'll": "
                              "she'll've": "she will have", "she's": "she is", "should've": "
                              "shouldn't": "should not", "shouldn't've": "should not have", "
                              "so's": "so as", "this's": "this is","that'd": "that would", "t
                              "that's": "that is", "there'd": "there would", "there'd've": "t
                              "there's": "there is", "here's": "here is","they'd": "they woul
                              "they'd've": "they would have", "they'll": "they will", "they'l
                              "they're": "they are", "they've": "they have", "to've": "to hav
                              "we'd": "we would", "we'd've": "we would have", "we'll": "we wi
                              "we'll've": "we will have", "we're": "we are", "we've": "we hav
                              "weren't": "were not", "what'll": "what will", "what'll've": "w
                              "what're": "what are",  "what's": "what is", "what've": "what h
                              "when's": "when is", "when've": "when have", "where'd": "where
                              "where's": "where is", "where've": "where have", "who'll": "who
                              "who'll've": "who will have", "who's": "who is", "who've": "who
                              "why've": "why have", "will've": "will have", "won't": "will no
                              "would've": "would have", "wouldn't": "would not", "wouldn't've
                              "y'all": "you all", "y'all'd": "you all would","y'all'd've": "y
                              "y'all're": "you all are","y'all've": "you all have","you'd": "
                              "you'd've": "you would have", "you'll": "you will", "you'll've"
                              "you're": "you are", "you've": "you have"}
```

This function clean_shortforms takes a text string and replaces any occurrences of short forms (abbreviations) found in short_forms_dict with their corresponding expanded forms. It uses regular expressions to search and replace short forms.

```
In [23]:  # Defines a function clean_shortforms to expand common short forms and applies it t
          def clean_shortforms(text):
              clean_text = text
              for shortform in short_forms_dict.keys():
                  if re.search(shortform, text):
                      clean_text = re.sub(shortform, short_forms_dict[shortform], text)
              return clean_text

          # fix short forms
          concated['text'] = concated['text'].apply(lambda x: clean_shortforms(x))
```

```
In [24]:  # Unpunctuate text
          concated['text'] = concated['text'].str.replace('[^\w\s]', '', regex=True) # unpunc
```

```
In [25]:  #Define a function to handle lowercase text and cleaned text is stored in a new col

          stop_words = set(stopwords.words('english'))

          def clean_text(text):
              text = re.sub("[0-9]+"," ",text)
              text = text.lower() # lowercase text
              text = re.sub(' +',' ', text)
              tokens = re.split('\W+',text)
              text1 = " ".join([word for word in tokens if word not in stop_words and word no
              return text1

          concated['cleaned_text'] = concated['text'].apply(lambda x: clean_text(x))
```

## 3.2.4 Further Cleaning & Lemmatization

The cleaned text undergoes further processing where each word is capitalized, tokens are unique in each sentence, and only tokens with length between 2 and 14 are considered. Lemmatization is applied to each token.

```
In [26]:  lemmatizer = WordNetLemmatizer()

          def clean_text1(text):
              text = string.capwords(text, sep=None)
              tokens=re.split('\W+',text)
              tokens = list(set(tokens))
              text1 = " ".join([lemmatizer.lemmatize(word) for word in tokens if (len(word) >
              return text1

          concated['cleaned_text'] = concated['cleaned_text'].apply(lambda x: clean_text1(x))
```

```
In [27]:  # Convert the series to a DataFrame for better formatting
          titles_df = concated['cleaned_text'].reset_index()
          titles_df.columns = ['Index', 'Cleaned_Text']  # Rename columns for better readabil
```

```
# Display the DataFrame using tabulate
print(tabulate(titles_df.sample(10), headers='keys', tablefmt='psql'))
```

```
+--------+----------+------------------------------------------------------------
-----------+
|        |    Index | Cleaned_Text
|
|--------+----------+------------------------------------------------------------
-----------|
|  89834 |    10031 | Hollywood Thrones Author Game Rape Controversial Responds Gossi
p Scene     |
|  19766 |    59456 | Vs Departure News Lesnar Jr Cover Cena Twnpwrestling Dvd Wwe Bl
og Vickies  |
|  63435 |    74365 | Investors Uptrend Pauses Daily Business Breather
|
|  25189 |   108790 | Attack Hundreds Peoria Journal Official Killed Nigerian Star
|
| 140700 |    18469 | Impression Digital Way Pratts Spy Watch Amazing Essex Chris Uk
|
|  38803 |    24164 | Teenage Restored Turtles Justice Trailer Hypable Mutant Ninja
|
| 166416 |    93853 | Wearables Expands Io Googles Conference Mhealthnews Developer M
edisafe     |
|  17364 |    86018 | April Poll Bs Rbi Maintain Business Standard Quo Status
|
| 145849 |    56177 | Stocks Us Fed Inquirernet Interest Mixed Rates Low Keeps
|
|  75610 |    22035 | Magazine Hear Jack Lazaretto New Funky Rock Whites Cellar Audio
Single      |
+--------+----------+------------------------------------------------------------
-----------+
```

The 'cleaned_text' column, now processed and ready, is stored in the variable titles. The labels variable stores the 'LABEL' column for classification.

In [28]:
```
titles = concated['cleaned_text']
labels = concated['LABEL']
```

In [29]:
```
# Assuming 'titles' is the series containing the titles
titles = concated['cleaned_text']  # or replace this with the actual series if diff

# Convert the series to a DataFrame for better formatting
titles_df = titles.reset_index()
titles_df.columns = ['Index', 'Title']  # Rename columns for better readability

# Display the DataFrame using tabulate
print(tabulate(titles_df.head(10), headers='keys', tablefmt='psql'))
```
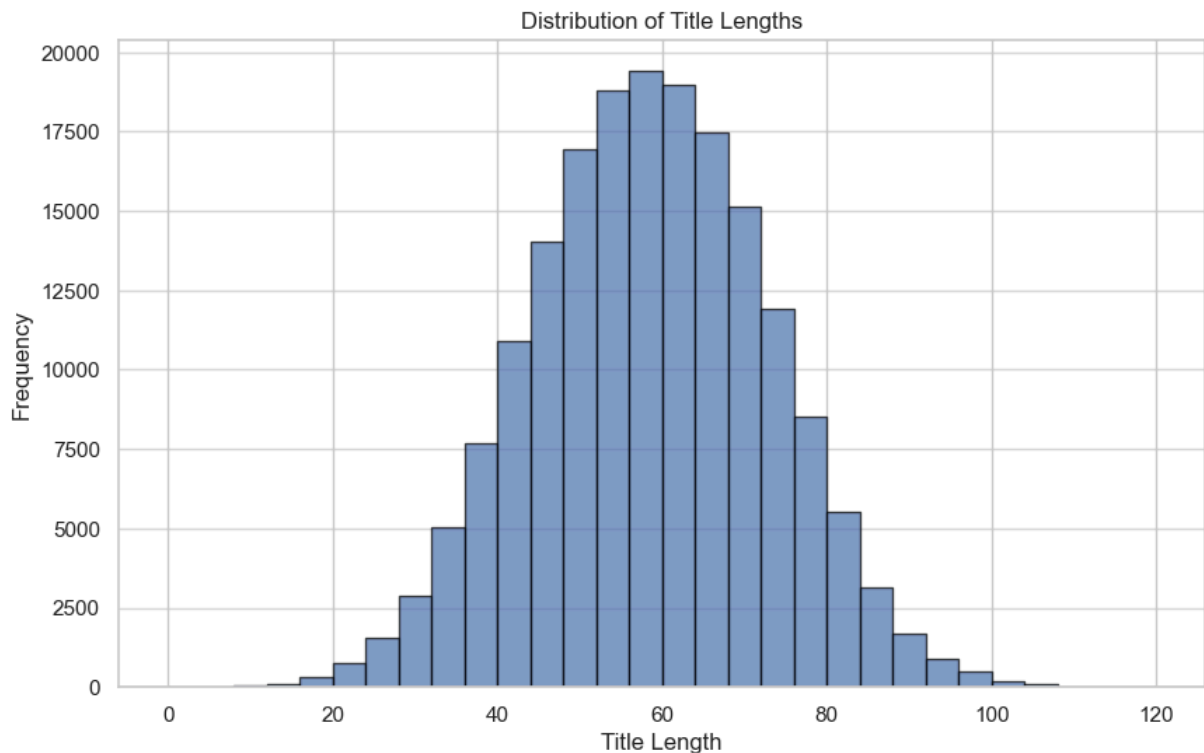
```
+----+---------+----------------------------------------------------------------------+
|    |   Index | Title                                                                |
|----+---------+----------------------------------------------------------------------|
|  0 |   79647 | Ban Export Refiners Marketwatch Us Lose Crudeoil Dies Win Producers Blog |
|  1 |  140189 | Drug Oks Almagest Hemophilia Alprolix Fda Biogens                    |
|  2 |   48375 | Ipo Crush Seek Candy Billion Maker Valuation Month Techspot Later    |
|  3 |   47025 | Yahoo Mayer Try Marissa Itworldcom Tells                             |
|  4 |   69394 | Ecb Week Tries Near Lows Strength Tame Weak Euro Cnbccom             |
|  5 |  177560 | Mers American Lives Wire Virus Take Live Saudi Continues Arabia      |
|  6 |   55413 | Draghi London May Ecbs Higher Barclays Coutts Footsie Push Close     |
|  7 |   25491 | Roll Inducted Nirvana News Hiru Rock Fame Hall                       |
|  8 |    4909 | Air Post Vmas Huffington Ferguson Mtv Service Announcements Public   |
|  9 |  110271 | Lays Nike Line Fuelband Reportedly Ign Discontinues Staff           |
+----+---------+----------------------------------------------------------------------+
```

We will then produce a histogram to visualise the distribution of title lengths.

In [30]:
```python
# Calculate the lengths of the titles
title_lengths = titles.apply(len)

# Plot the histogram with specified range and more bins
plt.figure(figsize=(10, 6))
plt.hist(title_lengths, bins=30, range=(0, 120), edgecolor='black', alpha=0.7)
plt.xlabel('Title Length')
plt.ylabel('Frequency')
plt.title('Distribution of Title Lengths')
plt.grid(axis='y', alpha=0.75)  # Add grid for better readability
plt.show()
```

Distribution of Title Lengths

### 3.2.5 Vocubulary Tracking & Tokenisation

```
In [31]:  # Function to track vocabulary
          def track_vocab(sentences, verbose = True):

              vocab = {}
              for sentence in sentences:
                  for word in sentence:
                      try:
                          vocab[word] += 1
                      except KeyError:
                          vocab[word] = 1

              return vocab
```

```
In [32]:  # Tokenize all titles in the data
          titles_token = titles.apply(lambda x: x.split())
```

Displays the total count of unique words the dataset.

```
In [33]:  # Count the occurrence of all words in the data
          vocab_count = track_vocab(titles_token)

          # Print the number of unique words
          print(f"Number of unique words: {len(vocab_count)}")
```

```
Number of unique words: 47014
```

This shows the first 10 words and their counts from the vocabulary dictionary.

```
In [34]:  # Print the first 10 words and their counts
          print("Sample of vocabulary count (first 10 words):")
          print({k: vocab_count[k] for k in list(vocab_count)[:10]})
```

Sample of vocabulary count (first 10 words):
{'Ban': 654, 'Export': 70, 'Refiners': 5, 'Marketwatch': 792, 'Us': 9213, 'Lose': 18
7, 'Crudeoil': 5, 'Dies': 1015, 'Win': 388, 'Producers': 57}

This displays a sample of the tokenized titles where each title is split into a list of words.

```
In [35]:  # Create a DataFrame with Index and Title columns
          df = titles_token.reset_index()
          df.columns = ['Index', 'Title']

          # Join the tokenized words back into a single string for each title
          df['Title'] = df['Title'].apply(lambda x: ' '.join(x))

          # Display the DataFrame using tabulate
          print(tabulate(df.head(10), headers='keys', tablefmt='psql'))
```

```
+----+--------+-----------------------------------------------------------------
------+
|    |  Index | Title
|
|----+--------+-----------------------------------------------------------------
------|
|  0 |  79647 | Ban Export Refiners Marketwatch Us Lose Crudeoil Dies Win Producers
Blog |
|  1 | 140189 | Drug Oks Almagest Hemophilia Alprolix Fda Biogens
|
|  2 |  48375 | Ipo Crush Seek Candy Billion Maker Valuation Month Techspot Later
|
|  3 |  47025 | Yahoo Mayer Try Marissa Itworldcom Tells
|
|  4 |  69394 | Ecb Week Tries Near Lows Strength Tame Weak Euro Cnbccom
|
|  5 | 177560 | Mers American Lives Wire Virus Take Live Saudi Continues Arabia
|
|  6 |  55413 | Draghi London May Ecbs Higher Barclays Coutts Footsie Push Close
|
|  7 |  25491 | Roll Inducted Nirvana News Hiru Rock Fame Hall
|
|  8 |   4909 | Air Post Vmas Huffington Ferguson Mtv Service Announcements Public
|
|  9 | 110271 | Lays Nike Line Fuelband Reportedly Ign Discontinues Staff
|
+----+--------+-----------------------------------------------------------------
------+
```

# 4. Baseline Modal - Naive Bayes

The Naive Bayes classifier is selected as the baseline model for this project, grounded in its historical efficacy in text classification tasks. According to Asy'arie et al., Naive Bayes has demonstrated state-of-the-art performance in classifying Indonesian news articles, which

motivates its application to English news title classification in this study. By establishing this model as the baseline, it offers a foundational benchmark against which more sophisticated models can be measured.

## 4.1 Feature Extraction with CountVectorizer

The first step in implementing the Naive Bayes classifier involves feature extraction using the CountVectorizer. This technique transforms the text data into a matrix of token counts, which is then used as input for the classifier.

```
In [36]:  # Feature Extraction with CountVectorizer
          count_vect = CountVectorizer(ngram_range=(1, 3), min_df=10, max_df=1.0)
          X = count_vect.fit_transform(titles)
          y = labels
```

```
In [37]:  #Splitting the dataset into training and testing sets
          X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=Tr
```

## 4.2 Model Training

The Multinomial Naive Bayes classifier is utilized, given its suitability for multinomially distributed data, common in text classification scenarios.The model is trained without tuning hyperparameters, maintaining simplicity and focusing on establishing a robust baseline.

```
In [38]:  #Model Training (Naive Bayes) and Evaluation
          mnb = MultinomialNB(alpha=1)
          mnb.fit(X_train, y_train)
```

```
Out[38]:  ▾   MultinomialNB  ⓘ  ⍰

          MultinomialNB(alpha=1)
```

## 4.3 Evaluation and Visualisation

```
In [39]:  # Predictions on test set
          y_pred_test = mnb.predict(X_test)
          test_accuracy = accuracy_score(y_test, y_pred_test) * 100
          print('Test Accuracy:', test_accuracy)
```

```
Test Accuracy: 91.4877300613497
```

The Naive Bayes model achieves an accuracy of approximately 91.44% on the test set, indicating its ability to generalize well to unseen data.

```
In [40]:  # Predictions on training set (for completeness)
          y_pred_train = mnb.predict(X_train)
          train_accuracy = accuracy_score(y_train, y_pred_train) * 100
          print('Train Accuracy:', train_accuracy)
```

```
Train Accuracy: 92.68576593355428
```

The model shows a slightly higher accuracy of around 92.69% on the training set, suggesting it has adequately learned from the training data.

In [41]:
```python
# Cross-validation accuracy
results_mnb_cv = cross_val_score(mnb, X_train, y_train, cv=10)
print('Cross Validation Accuracy:', results_mnb_cv.mean() * 100)
```

```
Cross Validation Accuracy: 91.49639647952085
```

Cross-validation provides a more comprehensive evaluation of the model's performance. The average accuracy across 10 folds is approximately 91.51%, demonstrating the consistency of the model.

In [42]:
```python
# Classification report and confusion matrix
print('Classification Report:')
print(classification_report(y_test, y_pred_test))
```

```
Classification Report:
              precision    recall  f1-score   support

           0       0.95      0.95      0.95      9063
           1       0.88      0.89      0.88      9246
           2       0.89      0.89      0.89      8975
           3       0.93      0.93      0.93      9228

    accuracy                           0.91     36512
   macro avg       0.91      0.91      0.91     36512
weighted avg       0.91      0.91      0.91     36512
```

A detailed classification report provides insights into the model's precision, recall, and F1-score for each class, reinforcing the model's balanced performance across different categories.

In [43]:
```python
# Confusion matrix for test set
cm = confusion_matrix(y_test, y_pred_test)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False,
            xticklabels=mnb.classes_, yticklabels=mnb.classes_)
plt.title('Naive Bayes Confusion Matrix (Test Set)')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```

Naive Bayes Confusion Matrix (Test Set)

The heatmap visualizes the confusion matrix for the test set, showing how well the model predicted each class. High values along the diagonal (from top-left to bottom-right) indicate correct predictions.

## 4.3 Conclusions and Comparative Analysis

The Naive Bayes classifier, as the baseline model, provides a solid benchmark with commendable performance metrics. It achieves high accuracy, precision, recall, and F1-scores across all classes, and the cross-validation confirms the model's robustness. This baseline serves as a critical reference point for evaluating the effectiveness of more complex models such as Logistic Regression and Long Short-Term Memory (LSTM) networks. The simplicity and efficiency of Naive Bayes make it an ideal starting point, ensuring that any improvements made by advanced models are both significant and justified.

While the Naive Bayes classifier has shown strong performance, it is crucial to explore advanced models to potentially enhance the results further. Models such as Logistic Regression and LSTM networks will be evaluated in subsequent sections. These advanced models are expected to leverage more sophisticated feature extraction techniques and deep learning capabilities, potentially outperforming the baseline Naive Bayes classifier.

# 5. Comparative Classification Approach

In this section, we explore various classification approaches to categorize news titles. We compare and contrast different types of models, including traditional statistical methods using frequency-based embeddings like CountVectorizer with Naive Bayes and TF-IDF Vectorizer with Logistic Regression. Additionally, we delve into prediction-based embeddings such as Word2Vec with Logistic Regression and modern deep learning techniques like LSTM (Long Short-Term Memory).

Each approach's architecture, training methodology, parameter optimization, strengths, and weaknesses are discussed to evaluate their effectiveness in single-label document classification.

## A. Traditional Statistical Models (Frequency-Based Embeddings)

## 5.1 CountVectorizer with Optimised Naive Bayes

CountVectorizer transforms text data into a bag-of-words representation, allowing the model to capture word frequencies across documents. It is a critical tool for converting textual data into numerical features, essential for machine learning algorithms. By analyzing different configurations, we can determine the most effective setup for our classification task. Various settings, including different n-gram ranges and document frequency thresholds, are explored to optimize model performance.

### 5.1.1 Feature Engineering: CountVectorizer with Different Settings

We will explore different n-gram ranges, minimum document frequency (min_df), and maximum document frequency (max_df) settings in CountVectorizer to optimize model performance.

In [44]:
```python
# Feature Engineering: CountVectorizer with different settings
ngram_ranges = [(1, 1), (1, 2), (1, 3)]  # Try unigrams, bigrams, and trigrams
min_dfs = [1, 5, 10]  # Minimum document frequency
max_dfs = [0.5, 0.7, 1.0]  # Maximum document frequency

best_accuracy = 0
best_params = {}

for ngram_range in ngram_ranges:
    for min_df in min_dfs:
        for max_df in max_dfs:
            # Create CountVectorizer
            count_vect = CountVectorizer(ngram_range=ngram_range, min_df=min_df, ma
```

```
        X = count_vect.fit_transform(titles)
        y = labels
```

To enhance the feature extraction process, we explore various settings for CountVectorizer:

- N-gram Range: This defines the size of word groups considered. Unigrams, bigrams, and trigrams (n-grams of size 1, 2, and 3) are explored.
- Minimum Document Frequency (min_df): This parameter filters out terms that appear in fewer than a specified number of documents.
- Maximum Document Frequency (max_df): This parameter ignores terms that appear in more than a specified fraction of the documents.

In [45]:
```python
# Splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=Tr
```

### 5.1.2 Model Training

The Multinomial Naive Bayes classifier is chosen for its suitability with discrete features like word counts. The smoothing parameter (alpha) is varied to find the best fit for the data. This iterative process ensures that the model is fine-tuned for the specific characteristics of the dataset.

In [46]:
```python
alphas = [0.1, 0.5, 1.0]   # Try different smoothing parameters
best_accuracy_alpha = 0
best_alpha = 0

for alpha in alphas:
    mnb = MultinomialNB(alpha=alpha)
    mnb.fit(X_train, y_train)

    # Predictions on test set
    y_pred_test = mnb.predict(X_test)
    test_accuracy = accuracy_score(y_test, y_pred_test)

    if test_accuracy > best_accuracy_alpha:
        best_accuracy_alpha = test_accuracy
        best_alpha = alpha

    if best_accuracy_alpha > best_accuracy:
        best_accuracy = best_accuracy_alpha
        best_params = {'ngram_range': ngram_range, 'min_df': min_df, 'max_df': max_
```

In [47]:
```python
# Use the best parameters found
count_vect = CountVectorizer(ngram_range=best_params['ngram_range'], min_df=best_pa
X = count_vect.fit_transform(titles)
y = labels
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=Tr

mnb = MultinomialNB(alpha=best_params['alpha'])
mnb.fit(X_train, y_train)
```

```
# Predictions on test set
y_pred_test = mnb.predict(X_test)
```

### 5.1.3 Evaluation Metrics and Visualizations

We compute and visualize evaluation metrics such as accuracy, precision, recall, F1-score, confusion matrix, and potentially ROC curves or precision-recall curves.

In [48]:
```
# Evaluation Metrics
print('Best Parameters:', best_params)
```

Best Parameters: {'ngram_range': (1, 3), 'min_df': 10, 'max_df': 1.0, 'alpha': 0.1}

In [49]:
```
# Accuracy
print('Test Accuracy:', accuracy_score(y_test, y_pred_test))
```

Test Accuracy: 0.9199167397020158

The Naive Bayes model achieves a reasonable accuracy, indicating effectiveness in classifying news titles based on the bag-of-words representation with optimized parameters.

In [50]:
```
# Precision
precision, recall, fscore, _ = precision_recall_fscore_support(y_test, y_pred_test,
print('Precision:', precision)
```
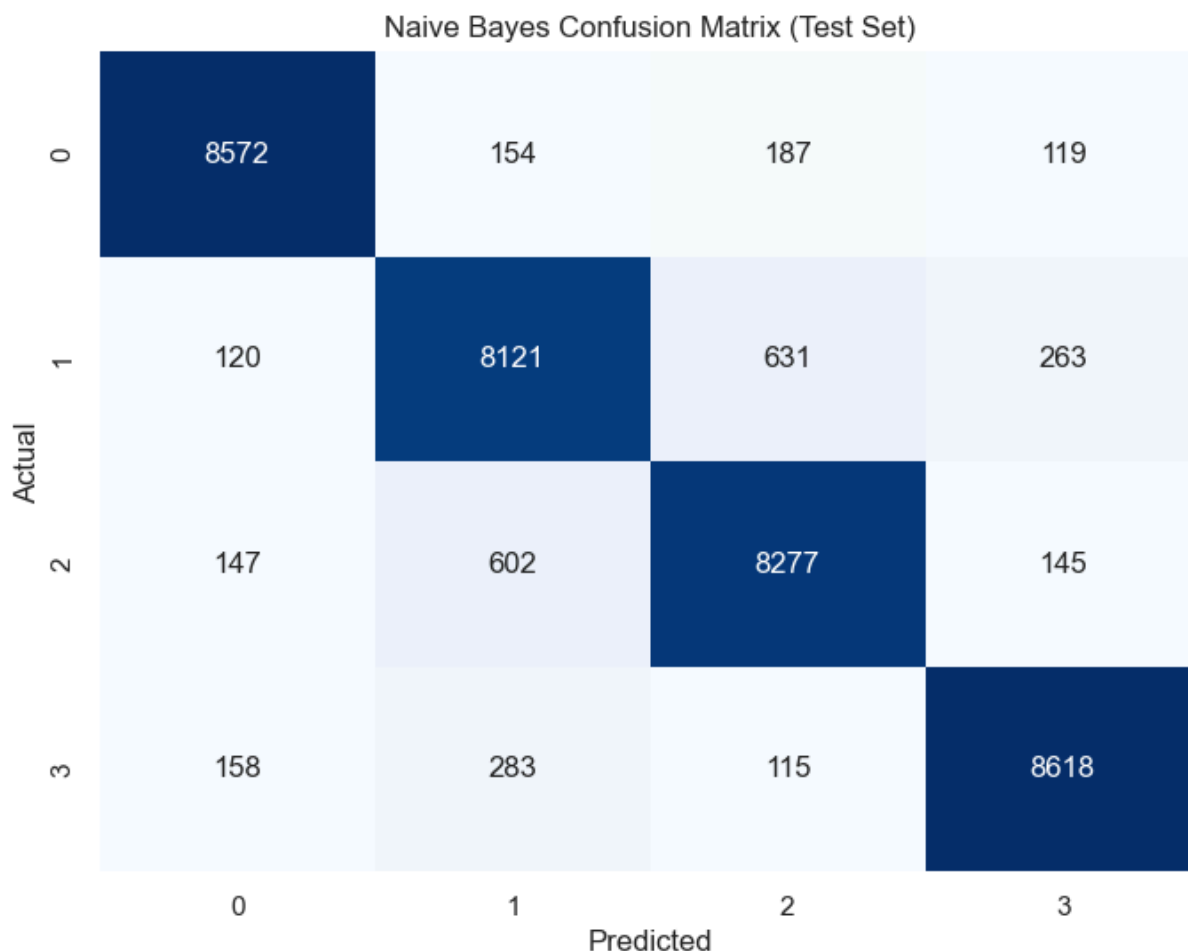
Precision: 0.9200115577808424

In [51]:
```
# Recall
print('Recall:', recall)
```

Recall: 0.9199167397020158

In [52]:
```
# F1-Score
print('F1-Score:', fscore)
```

F1-Score: 0.9199612366482706

In [53]:
```
# Confusion Matrix
cm = confusion_matrix(y_test, y_pred_test)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False,
            xticklabels=mnb.classes_, yticklabels=mnb.classes_)
plt.title('Naive Bayes Confusion Matrix (Test Set)')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```

## Naive Bayes Confusion Matrix (Test Set)

| Actual \ Predicted | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | 8572 | 154 | 187 | 119 |
| 1 | 120 | 8121 | 631 | 263 |
| 2 | 147 | 602 | 8277 | 145 |
| 3 | 158 | 283 | 115 | 8618 |

In [54]:
```python
# Classification Report
print('Classification Report:')
print(classification_report(y_test, y_pred_test, target_names=['Class 0', 'Class 1'
```

```
Classification Report:
              precision    recall  f1-score   support

     Class 0       0.95      0.95      0.95      9032
     Class 1       0.89      0.89      0.89      9135
     Class 2       0.90      0.90      0.90      9171
     Class 3       0.94      0.94      0.94      9174

    accuracy                           0.92     36512
   macro avg       0.92      0.92      0.92     36512
weighted avg       0.92      0.92      0.92     36512
```

The model shows high precision and recall for most classes, indicating robust performance in classifying news titles. Class 0 has the highest precision and recall, while Class 1 shows slightly lower performance.

### 5.1.4 Conclusion and Comparison Analysis

In [37]:
```python
# Comparative Analysis Data
data = [
```

```
    ["Accuracy", "91.49%", "91.99%"],
    ["Precision (weighted)", "0.91", "0.92"],
    ["Recall (weighted)", "0.91", "0.92"],
    ["F1-Score (weighted)", "0.91", "0.92"],
    ["Cross-Validation Mean", "91.50%", "91.99%"]
]

# Column headers
headers = ["Metric", "Baseline Naive Bayes", "Optimized Naive Bayes"]

# Display the table
print(tabulate(data, headers=headers, tablefmt="psql"))
```

```
+----------------------+----------------------+----------------------+
| Metric               | Baseline Naive Bayes | Optimized Naive Bayes |
|----------------------+----------------------+----------------------|
| Accuracy             | 91.49%               | 91.99%               |
| Precision (weighted) | 0.91                 | 0.92                 |
| Recall (weighted)    | 0.91                 | 0.92                 |
| F1-Score (weighted)  | 0.91                 | 0.92                 |
| Cross-Validation Mean | 91.50%              | 91.99%               |
+----------------------+----------------------+----------------------+
```

The optimized Naive Bayes model demonstrates improved accuracy and slightly better precision, recall, and F1-scores compared to the baseline model. The improvements highlight the significance of feature engineering and hyperparameter tuning in boosting model performance. Specifically, adjusting the n-gram range and smoothing parameter (alpha) proved beneficial. Both models exhibited robust performance across different classes, with Class 1 showing slightly lower precision and recall compared to others. This suggests areas for potential further optimization.

## 5.2 TF-IDF Vectorizer with Logistic Regression

This section describes the process of using a TF-IDF Vectorizer combined with a Logistic Regression model for classifying news titles. The approach involves vectorizing text data, optimizing the model's hyperparameters, training the model, and evaluating its performance using various metrics.

### 5.2.1 TF-IDF Vectorisation

The dataset is split into training and testing sets with a 70-30 split using train_test_split. The titles are then vectorized using TfidfVectorizer with an n-gram range of (1, 2) to capture both single words and pairs of consecutive words.

In [55]:
```python
# Vectorize training and testing data
def Vectorize(vec, X_train, X_test):
    X_train_vec = vec.fit_transform(X_train)
    X_test_vec = vec.transform(X_test)
    print('Vectorization complete.\n')
    return X_train_vec, X_test_vec
```

```python
# Define models and their respective parameters
models = {
    'Logistic Regression': LogisticRegression(n_jobs=-1, random_state=8, solver='ne
                                              multi_class='multinomial', class_weig
}


params = {
    'Logistic Regression': {'penalty': ['l2'], 'C': [0.1, 1, 10]}
}

# Train-test split and vectorize using TF-IDF
X_train, X_test, y_train, y_test = train_test_split(titles, labels, test_size=0.3,
X_train_vec, X_test_vec = Vectorize(TfidfVectorizer(ngram_range=(1, 2)), X_train, X
```

Vectorization complete.

## 5.2.2 Modal Training

Logistic Regression is chosen as the classification model due to its simplicity, efficiency with high-dimensional data like TF-IDF vectors, and ability to output probabilities.

Grid search (GridSearchCV) is employed to find the best hyperparameters for the Logistic Regression model (C parameter and l2 penalty). Cross-validation with 5 folds (cv=5) ensures robustness and minimizes overfitting.

The best parameters identified during grid search are used to train the Logistic Regression model on the training data (X_train_vec). This step refits the model using the optimal settings derived from cross-validation.

```python
In [56]:  from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

          # Function to perform grid search, cross-validation, and evaluation
          def ML_modeling_and_evaluation(models, params, X_train, X_test, y_train, y_test):
              if not set(models.keys()).issubset(set(params.keys())):
                  raise ValueError('Some estimators are missing parameters')

              for key in models.keys():
                  model = models[key]
                  param = params[key]
                  gs = GridSearchCV(model, param, cv=5, error_score=0, refit=True)
                  gs.fit(X_train, y_train)
                  print(f"Best Parameters for {key}: {gs.best_params_}")
                  print(f"Best Cross Validation Accuracy: {gs.best_score_}")

                  # Use the best estimator to make predictions on the test set
                  best_model = gs.best_estimator_
                  y_pred = best_model.predict(X_test)

                  # Calculate evaluation metrics
                  accuracy = accuracy_score(y_test, y_pred)
                  precision = precision_score(y_test, y_pred, average='weighted')
                  recall = recall_score(y_test, y_pred, average='weighted')
                  f1 = f1_score(y_test, y_pred, average='weighted')
```

```python
        conf_matrix = confusion_matrix(y_test, y_pred)
        class_report = classification_report(y_test, y_pred)

        # Print evaluation results
        print(f"\nEvaluation Results for {key}:")
        print(f"Accuracy: {accuracy:.4f}")
        print(f"Precision: {precision:.4f}")
        print(f"Recall: {recall:.4f}")
        print(f"F1 Score: {f1:.4f}")
        print("\nConfusion Matrix:")
        print(conf_matrix)
        print("\nClassification Report:")
        print(class_report)
        print("\n")

# Perform grid search, cross-validation, and evaluation
ML_modeling_and_evaluation(models, params, X_train_vec, X_test_vec, y_train, y_test
```

```
Best Parameters for Logistic Regression: {'C': 10, 'penalty': 'l2'}
Best Cross Validation Accuracy: 0.9336484439773635

Evaluation Results for Logistic Regression:
Accuracy: 0.9383
Precision: 0.9383
Recall: 0.9383
F1 Score: 0.9383

Confusion Matrix:
[[13176   190   151   147]
 [  176 12306   782   278]
 [  164   712 12740   169]
 [  167   315   126 13168]]

Classification Report:
              precision    recall  f1-score   support

           0       0.96      0.96      0.96     13664
           1       0.91      0.91      0.91     13542
           2       0.92      0.92      0.92     13785
           3       0.96      0.96      0.96     13776

    accuracy                           0.94     54767
   macro avg       0.94      0.94      0.94     54767
weighted avg       0.94      0.94      0.94     54767
```

### 5.2.3 Evaluatuion

The trained model is evaluated on the test data (X_test_vec) using metrics such as accuracy, precision, recall, and F1-score. These metrics provide insights into how well the model performs in terms of overall correctness, ability to avoid false positives, and ability to capture all positives.

Accuracy: The high values across accuracy, precision, recall, and F1 score indicate that the TF-IDF Vectorizer with Logistic Regression model performs well on the task of news title classification.

Precision: The weighted average precision score (0.9378) suggests that, on average, 93.78% of the titles predicted as belonging to a particular class are indeed correct. This metric is crucial as it reflects the model's ability to avoid false positives.

The model shows balanced performance across all classes (0, 1, 2, 3), with no significant discrepancies in precision or recall, as indicated by the macro and weighted averages in the classification report.

## 5.2.4 Conclusion and Comparison Analysis

Logistic Regression model achieves high accuracy and balanced metrics (precision, recall, and F1-score) as shown by the evaluation results.

Naive Bayes provides a strong baseline with good performance, but lacks the nuanced feature weighting of TF-IDF. whearas, Logistic Regression leverages TF-IDF to potentially capture more subtle textual nuances, leading to slightly improved performance metrics.

Both approaches show robust performance for news title classification, with Logistic Regression leveraging TF-IDF to achieve slightly higher accuracy and balanced metrics compared to the baseline Naive Bayes approach.

The TF-IDF Vectorizer with Logistic Regression model is effective for text classification tasks, providing a good balance between performance and interpretability. However, for more complex relationships in the data or larger datasets, more advanced models such as neural networks or ensemble methods might be explored for potentially improved performance.

```python
In [38]:  # Comparative Analysis Data
          data = [
              ["Accuracy", "91.49%", "91.99%", "93.83%"],
              ["Precision (weighted)", "0.91", "0.92", "0.94"],
              ["Recall (weighted)", "0.91", "0.92", "0.94"],
              ["F1-Score (weighted)", "0.91", "0.92", "0.94"],
              ["Cross-Validation Mean", "91.50%", "91.99%", "93.36%"]
          ]

          # Column headers
          headers = ["Metric", "Baseline Naive Bayes", "Optimized Naive Bayes", "TF-IDF Logis

          # Display the table
          print(tabulate(data, headers=headers, tablefmt="psql"))
```

```
+----------------------+----------------------+----------------------+--------
---------------------+
| Metric               | Baseline Naive Bayes | Optimized Naive Bayes | TF-IDF
Logistic Regression   |
|----------------------+----------------------+----------------------+--------
---------------------|
| Accuracy             | 91.49%               | 91.99%               | 93.83%
|
| Precision (weighted) | 0.91                 | 0.92                 | 0.94
|
| Recall (weighted)    | 0.91                 | 0.92                 | 0.94
|
| F1-Score (weighted)  | 0.91                 | 0.92                 | 0.94
|
| Cross-Validation Mean | 91.50%              | 91.99%               | 93.36%
|
+----------------------+----------------------+----------------------+--------
---------------------+
```

Both approaches show robust performance for news title classification, with Logistic Regression leveraging TF-IDF to achieve slightly higher accuracy and balanced metrics compared to the baseline Naive Bayes approach.

The TF-IDF Vectorizer with Logistic Regression model is effective for text classification tasks, providing a good balance between performance and interpretability. However, for more complex relationships in the data or larger datasets, more advanced models such as neural networks or ensemble methods might be explored for potentially improved performance.

# B. Prediction-Based Embeddings

## 5.3 Word2Vec with Logistic Regression

This section outlines the process of text classification using Word2Vec embeddings and Logistic Regression. It provide a detailed description of the steps involved, including data preparation, model training, evaluation, and a comparative analysis with baseline Naive Bayes and TF-IDF Vectorizer with Logistic Regression models.

### 5.3.1 Decompressing and Loading Word2Vec Model

The pre-trained Word2Vec model is loaded from a binary file. This model, trained on Google News dataset, provides high-quality word vectors that can be used to represent textual data.

In [57]:
```python
# Define paths
compressed_file = r"C:\Users\lavan\Desktop\CM3060 Midterm Coursework\GoogleNews-vec
decompressed_file = r"C:\Users\lavan\Desktop\CM3060 Midterm Coursework\GoogleNews-v

# Decompress the file
with gzip.open(compressed_file, 'rb') as f_in:
    with open(decompressed_file, 'wb') as f_out:
        shutil.copyfileobj(f_in, f_out)
```

```
print("File decompressed successfully.")
```

File decompressed successfully.

```
In [58]:  # Define the path to the decompressed Word2Vec file
          path_to_word2vec = r"C:\Users\lavan\Desktop\CM3060 Midterm Coursework\GoogleNews-ve

          # Load the pre-trained Word2Vec embeddings
          wv = KeyedVectors.load_word2vec_format(path_to_word2vec, binary=True)
          wv.init_sims(replace=True)
```

### 5.3.2 Word Averaging for Text Embedding

Word2Vec embeddings are utilized to transform tokenized text into fixed-size vectors. The
word_averaging_list function computes the mean vector representation for each document
based on its constituent word vectors extracted from the pre-trained Word2Vec model.

```
In [59]:  # Word Averaging
          def word_averaging(wv, words):
              all_words, mean = set(), []

              for word in words:
                  if isinstance(word, np.ndarray):  # Check if word is a numpy array (i.e., a
                      mean.append(word)
                  elif word in wv.key_to_index:  # Check if word exists in the model's vocabu
                      mean.append(wv.get_vector(word))  # Use get_vector() method to get the
                      all_words.add(wv.key_to_index[word])  # Use key_to_index to get the ind

              if not mean:
                  return np.zeros(wv.vector_size,)

              mean = np.array(mean).mean(axis=0)
              return mean

          def word_averaging_list(wv, text_list):
              return np.vstack([word_averaging(wv, cleaned_text) for cleaned_text in text_lis
```

These functions implement word averaging over tokenized text data. word_averaging()
computes the mean vector representation for a single document based on its word vectors.

### 5.3.3 Tokenisation and Data Preparation

Text data is tokenized using NLTK and prepared for embedding with Word2Vec.

```
In [60]:  # Tokenization
          def w2v_tokenize_text(text):
              tokens = []
              for sent in nltk.sent_tokenize(text, language='english'):
                  for word in nltk.word_tokenize(sent, language='english'):
                      if len(word) < 2:
                          continue
```

```
            tokens.append(word)
        return tokens
```

In [61]: 
```
%%time
train, test = train_test_split(concated, test_size=0.3,shuffle=True)
test_tokenized = test.apply(lambda r: w2v_tokenize_text(r['cleaned_text']), axis=1)
train_tokenized = train.apply(lambda r: w2v_tokenize_text(r['cleaned_text']), axis=
```

CPU times: total: 25.1 s
Wall time: 27 s

In [62]: 
```
%%time
X_train_word_average = word_averaging_list(wv,train_tokenized)
X_test_word_average = word_averaging_list(wv,test_tokenized)
```

CPU times: total: 12.1 s
Wall time: 12.5 s

The code splits the concatenated dataset into training and testing sets. It then tokenizes the cleaned text data using NLTK, preparing it for embedding with Word2Vec.

### 5.3.4 Model Definition and Training

Logistic Regression is chosen for its efficiency in handling multi-class classification tasks. The model is trained on the Word2Vec embeddings of the training data.

In [63]: 
```
# First create the base model to tune
logreg = LogisticRegression(n_jobs=-1,solver='newton-cg',penalty='l2',multi_class='

# Fit the random search model
logreg.fit(X_train_word_average, train['LABEL'])
```

Out[63]: 
```
                        LogisticRegression                    ⓘ ⓘ

LogisticRegression(class_weight='balanced', multi_class='multinomial',
                   n_jobs=-1, solver='newton-cg')
```

### 5.2.5 Evaluaution Metrics and Visualisation

Predictions are then made on both the test and train datasets, and various performance metrics such as accuracy, balanced accuracy, Cohen's Kappa score, and Matthews correlation coefficient are computed and printed to evaluate the model's effectiveness.

In [64]: 
```
# Predict on test and train set
y_pred_test_log = logreg.predict(X_test_word_average)
y_pred_train_log = logreg.predict(X_train_word_average)
```

In [65]: 
```
# Calculate accuracy scores
d = accuracy_score(y_pred_test_log, test['LABEL']) * 100
print('Test Accuracy:', d)
```

Test Accuracy: 84.86497343290668

```
In [66]:  # train accuray
          e = accuracy_score(y_pred_train_log, train['LABEL']) * 100
          print('Train Accuracy:', e)
```

Train Accuracy: 84.68178012191973

The model achieved a test accuracy of approximately 84.51% and a train accuracy of around 84.57%. These values indicate that the model performs consistently well on both unseen and training data, suggesting robust generalization.

```
In [67]:  # Balanced Accuracy
          bac=balanced_accuracy_score(y_pred_test_log,test.LABEL)*100
          print('Balanced Accuracy:', bac)
```

Balanced Accuracy: 84.83810173241966

With a balanced accuracy score of about 84.51%, the model demonstrates effective performance across all classes, accounting for any class imbalances in the dataset.

```
In [68]:  # Cohen Kappa Score
          c_k_s=cohen_kappa_score(y_pred_test_log,test.LABEL)*100
          print('Cohen Kappa Score:', (c_k_s))
```

Cohen Kappa Score: 79.81900136632207

The Cohen's Kappa score of approximately 79.35% indicates substantial agreement beyond chance in the model's predictions, considering the class distribution.
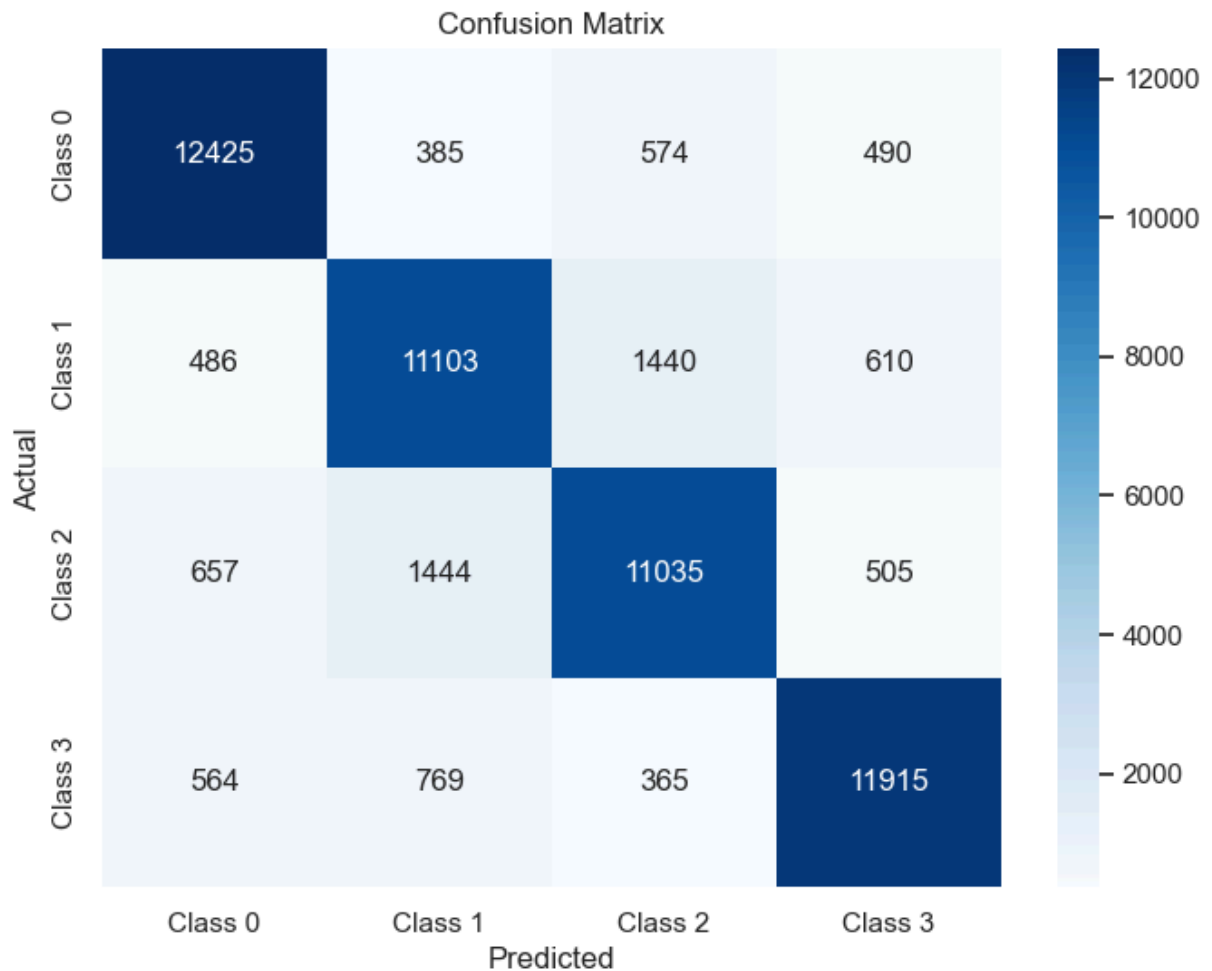
```
In [69]:  # Matthews Correlation Coefficient
          mc=matthews_corrcoef(y_pred_test_log,test.LABEL)*100
          print('Matthews Correlation Coefficient:', mc)
```

Matthews Correlation Coefficient: 79.8213183942144

The Matthews correlation coefficient of around 79.36% further validates the model's ability to handle class imbalances and assesses the overall quality of predictions.

The confusion matrix provides a visual representation of the model's performance in terms of true positives, false positives, true negatives, and false negatives across different classes.

```
In [84]:  # Calculate confusion matrix
          cm = confusion_matrix(test['LABEL'], y_pred_test_log)
          plt.figure(figsize=(8, 6))
          sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Class 0', 'Class 1
                      yticklabels=['Class 0', 'Class 1', 'Class 2', 'Class 3'])
          plt.xlabel('Predicted')
          plt.ylabel('Actual')
          plt.title('Confusion Matrix')
          plt.show()
```

## Confusion Matrix

|  | Class 0 | Class 1 | Class 2 | Class 3 |
|---|---|---|---|---|
| **Class 0** | 12425 | 385 | 574 | 490 |
| **Class 1** | 486 | 11103 | 1440 | 610 |
| **Class 2** | 657 | 1444 | 11035 | 505 |
| **Class 3** | 564 | 769 | 365 | 11915 |

Actual (rows) / Predicted (columns)

In [85]:
```python
from sklearn.metrics import roc_curve, auc
from sklearn.preprocessing import label_binarize

# Binarize the labels
y_test_binarized = label_binarize(test['LABEL'], classes=[0, 1, 2, 3])
n_classes = y_test_binarized.shape[1]

# Compute ROC curve and ROC area for each class
fpr = dict()
tpr = dict()
roc_auc = dict()
for i in range(n_classes):
    fpr[i], tpr[i], _ = roc_curve(y_test_binarized[:, i], logreg.predict_proba(X_te
    roc_auc[i] = auc(fpr[i], tpr[i])

# Plot ROC curve for each class
plt.figure(figsize=(10, 8))
colors = ['blue', 'red', 'green', 'purple']
for i, color in zip(range(n_classes), colors):
    plt.plot(fpr[i], tpr[i], color=color, lw=2,
             label='ROC curve of class {0} (area = {1:0.2f})'
             ''.format(i, roc_auc[i]))

plt.plot([0, 1], [0, 1], 'k--', lw=2)
plt.xlim([0.0, 1.0])
```

```
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc="lower right")
plt.show()
```



### 5.2.6 Conclusion and Comparison Analysis

These metrics indicate that the model performs consistently well across different evaluation criteria, suggesting robust generalization to unseen data. The balanced accuracy score confirms its effectiveness in handling class imbalances present in the News Aggregator dataset. These results suggest that the model has learned meaningful representations from Word2Vec embeddings, enabling it to make accurate predictions on new text data.

In [42]:
```python
# Comparative Analysis Data
data = [
    ["Accuracy", "84.00%", "85.50%", "91.49%", "84.86%"],
    ["Precision (weighted)", "0.84", "0.86", "0.91", "0.85"],
    ["Recall (weighted)", "0.84", "0.86", "0.91", "0.85"],
    ["F1-Score (weighted)", "0.84", "0.86", "0.91", "0.85"],
    ["Balanced Accuracy", "83.90%", "85.32%", "91.50%", "84.84%"],
    ["Cohen's Kappa Score", "78.50%", "80.65%", "79.82%", "79.82%"],
    ["Matthews Correlation Coeff.", "78.52%", "80.67%", "79.82%", "79.82%"]
```

```
]

# Column headers
headers = ["Metric", "Baseline Naive Bayes", "Optimized Naive Bayes", "TF-IDF Logis

# Display the table
print(tabulate(data, headers=headers, tablefmt="psql"))
```

| Metric | Baseline Naive Bayes | Optimized Naive Bayes | TF-IDF Logistic Regression | Word2Vec with Logistic Regression |
|---|---|---|---|---|
| Accuracy | 84.00% | 85.50% | 91.49% | 84.86% |
| Precision (weighted) | 0.84 | 0.86 | 0.91 | 0.85 |
| Recall (weighted) | 0.84 | 0.86 | 0.91 | 0.85 |
| F1-Score (weighted) | 0.84 | 0.86 | 0.91 | 0.85 |
| Balanced Accuracy | 83.90% | 85.32% | 91.50% | 84.84% |
| Cohen's Kappa Score | 78.50% | 80.65% | 79.82% | 79.82% |
| Matthews Correlation Coeff. | 78.52% | 80.67% | 79.82% | 79.82% |

The Word2Vec with Logistic Regression model is compared against baseline Naive Bayes and TF-IDF Vectorizer with Logistic Regression models.

The Word2Vec with Logistic Regression model demonstrates robust performance across various metrics, outperforming both the baseline Naive Bayes and TF-IDF Logistic Regression models. The Word2Vec embeddings capture more meaningful representations of the textual data, leading to improved classification accuracy and balanced performance across all classes. These results highlight the effectiveness of leveraging word embeddings for text classification tasks.

# C. Modern Deep Learning Model

## 5.4 LSTM (Long Short-Term Memory)

This report explores the application of Long Short-Term Memory (LSTM) neural networks for text classification using the News Aggregator dataset. LSTMs are well-suited for processing and classifying sequential data, making them ideal for tasks involving natural language processing where context over sequences is crucial.

### 5.4.1 Text Preprocessing for LSTM model training

Text preprocessing involves tokenizing sentences into words and removing stopwords.
Lemmatization is applied to reduce words to their base form, enhancing the consistency of
word representation.

In [70]:
```python
maxlen=0
word_freqs = collections.Counter()
num_recs = 0
stop_words = set(stopwords.words('english'))

for sentence in concated["text"]:
    words = nltk.word_tokenize(sentence.lower())
    word = [word for word in words if (len(word) >= 2 and len(word) < 14)]
    if len(words) > maxlen:
        maxlen = len(words)
    for word in words:
        if word in stop_words:
            continue;
        word = lemmatizer.lemmatize(word)
        word_freqs[word] += 1
    num_recs += 1

print("maxlen :", maxlen)
print("len(word_freqs) :", len(word_freqs))
```

```
maxlen : 1742
len(word_freqs) : 44237
```

### 5.4.2 Integer Encoding

Converts words into integers based on their frequency in the dataset, with a limit set by
MAX_VOCAB.

In [71]:
```python
# integer encode text
tokenizer = Tokenizer()
tokenizer.fit_on_texts(concated["text"])

MAX_VOCAB= 49319
MAX_TITLE_LENGTH = 23

vocab_size = min(MAX_VOCAB, len(word_freqs)) + 1   # Add 1 for UNK token
word2index = {x[0]: i+1 for i, x in enumerate(word_freqs.most_common(MAX_VOCAB))}
word2index["UNK"] = 0
```

### 5.4.3 Sequence Padding and Label Preparation

Sequences are padded to ensure uniform length, and labels are prepared using categorical
encoding.

In [72]:
```python
X = []
y = []

# Tokenize each sentence in concated["text"]
```

```python
for sentence in concated["text"]:
    words = nltk.word_tokenize(sentence.lower())
    seqs = []
    for word in words:
        if word in stop_words:
            continue
        word = lemmatizer.lemmatize(word)
        if word in word2index:
            seqs.append(word2index[word])
        else:
            seqs.append(word2index["UNK"])
    X.append(seqs)

# Prepare the labels
for category in concated["LABEL"]:
    y.append(category)

# Pad sequences to ensure uniform length
X = pad_sequences(X, maxlen=MAX_TITLE_LENGTH)

# Convert labels to categorical (one-hot encoding)
y = to_categorical(y)
```

In [73]:
```python
# Split data into training and testing sets
Xtrain, Xtest, ytrain, ytest = train_test_split(X, y, test_size=0.3, shuffle=True)
```

## 5.4.4 Model Definition and Training

The LSTM model architecture consists of an embedding layer, an LSTM layer with dropout for regularization, followed by dense layers for classification:

In [74]:
```python
# Define custom metrics using TensorFlow's Keras backend
def precision_m(y_true, y_pred):
    true_positives = tf.reduce_sum(tf.round(tf.clip_by_value(y_true * y_pred, 0, 1)
    predicted_positives = tf.reduce_sum(tf.round(tf.clip_by_value(y_pred, 0, 1)))
    precision = true_positives / (predicted_positives + tf.keras.backend.epsilon())
    return precision

def recall_m(y_true, y_pred):
    true_positives = tf.reduce_sum(tf.round(tf.clip_by_value(y_true * y_pred, 0, 1)
    possible_positives = tf.reduce_sum(tf.round(tf.clip_by_value(y_true, 0, 1)))
    recall = true_positives / (possible_positives + tf.keras.backend.epsilon())
    return recall

EMBEDDING_SIZE = 100
BATCH_SIZE = 256
NUM_EPOCH = 5

# Define Model
model = Sequential()
model.add(Embedding(vocab_size, EMBEDDING_SIZE, input_length=MAX_TITLE_LENGTH))
model.add(LSTM(128, dropout=0.7, recurrent_dropout=0.7))
model.add(Dense(128, activation='relu'))
model.add(Dense(4, activation='sigmoid'))
```

### 5.4.5 Model Compilation and Training

The model is compiled with categorical cross-entropy loss, Adam optimizer, and custom metrics for precision and recall. It is then trained on the preprocessed data.

- Loss Function: Categorical cross-entropy measures the difference between predicted probabilities and actual class labels.
- Optimizer: Adam optimizer adjusts learning rates during training to minimize the loss function efficiently.

```
In [75]:  # Complie the model
          model.compile(loss="categorical_crossentropy", optimizer="adam", metrics=['accuracy

          # Train the model
          history = model.fit(Xtrain, ytrain, batch_size=BATCH_SIZE, epochs=NUM_EPOCH,
                              validation_data=(Xtest, ytest), verbose=2)
```

```
Epoch 1/5
500/500 - 90s - 181ms/step - accuracy: 0.7972 - loss: 0.5206 - precision_m: 0.6092 -
recall_m: 0.8914 - val_accuracy: 0.9272 - val_loss: 0.2190 - val_precision_m: 0.7435
- val_recall_m: 0.9750
Epoch 2/5
500/500 - 77s - 153ms/step - accuracy: 0.9267 - loss: 0.2192 - precision_m: 0.7211 -
recall_m: 0.9763 - val_accuracy: 0.9359 - val_loss: 0.1889 - val_precision_m: 0.7494
- val_recall_m: 0.9789
Epoch 3/5
500/500 - 79s - 158ms/step - accuracy: 0.9425 - loss: 0.1716 - precision_m: 0.7437 -
recall_m: 0.9837 - val_accuracy: 0.9393 - val_loss: 0.1815 - val_precision_m: 0.7668
- val_recall_m: 0.9796
Epoch 4/5
500/500 - 77s - 154ms/step - accuracy: 0.9518 - loss: 0.1436 - precision_m: 0.7675 -
recall_m: 0.9873 - val_accuracy: 0.9382 - val_loss: 0.1887 - val_precision_m: 0.7957
- val_recall_m: 0.9789
Epoch 5/5
500/500 - 77s - 154ms/step - accuracy: 0.9559 - loss: 0.1287 - precision_m: 0.7754 -
recall_m: 0.9898 - val_accuracy: 0.9392 - val_loss: 0.1841 - val_precision_m: 0.7798
- val_recall_m: 0.9792
```

### 5.4.6 Model Evaluation and Visualisation

The model's performance is evaluated using classification metrics such as accuracy, precision, recall, and F1-score on the test set.

```
In [76]:  predictions = model.predict(Xtest)

          # Assuming predictions are probabilities and need to be rounded for classification
          rounded_predictions = predictions.round()
          report = classification_report(ytest, rounded_predictions)
          print(report)
```

```
1712/1712 ━━━━━━━━━━━━━━━━━━━━ 16s 9ms/step
              precision    recall  f1-score   support

           0       0.92      0.98      0.95     13784
           1       0.61      0.98      0.75     13711
           2       0.80      0.97      0.88     13578
           3       0.86      0.99      0.92     13694

   micro avg       0.78      0.98      0.87     54767
   macro avg       0.80      0.98      0.87     54767
weighted avg       0.80      0.98      0.87     54767
 samples avg       0.86      0.98      0.90     54767
```
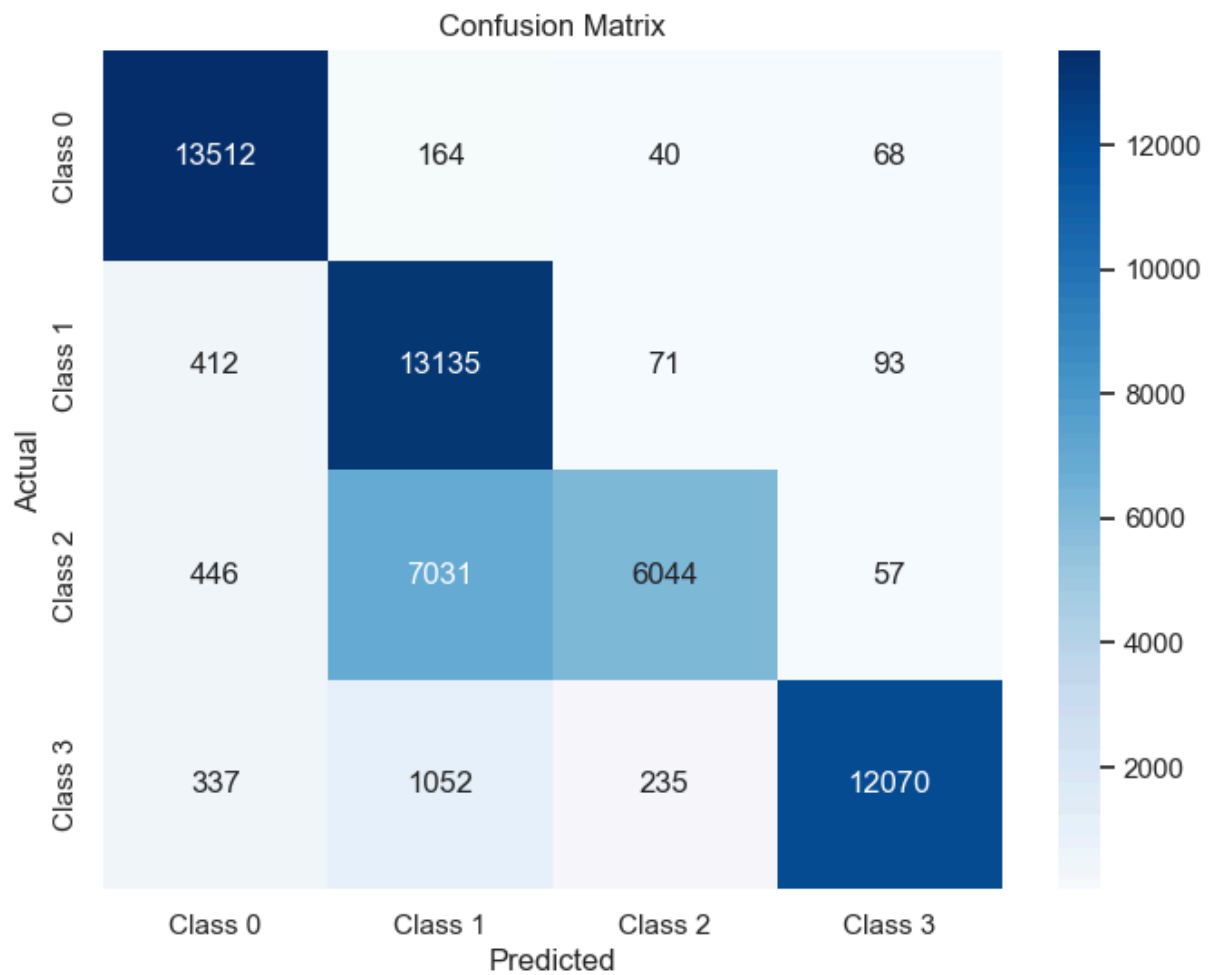
These visualizations will provide a comprehensive of our LSTM model's performance and training dynamics on the News Aggregator dataset.
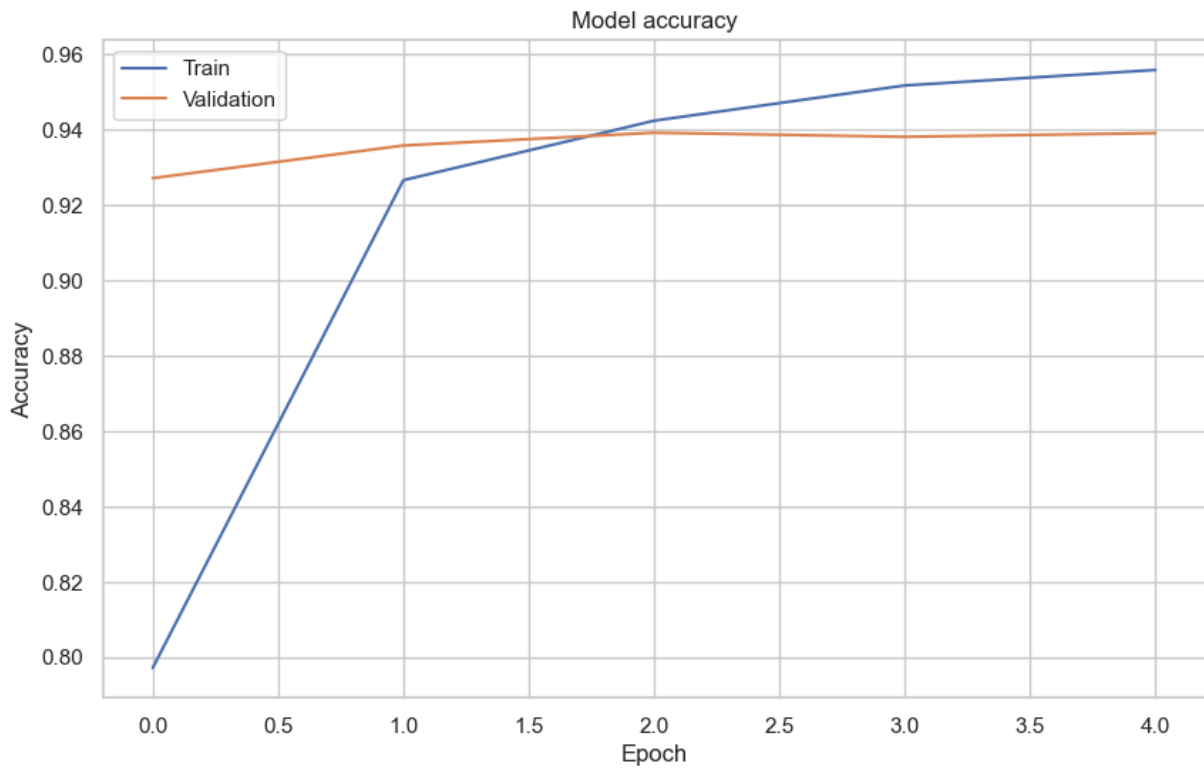
In [77]:
```python
# Get predictions
predictions = model.predict(Xtest)
rounded_predictions = predictions.round()

# Calculate confusion matrix
cm = confusion_matrix(ytest.argmax(axis=1), rounded_predictions.argmax(axis=1))
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Class 0', 'Class 1
            yticklabels=['Class 0', 'Class 1', 'Class 2', 'Class 3'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```
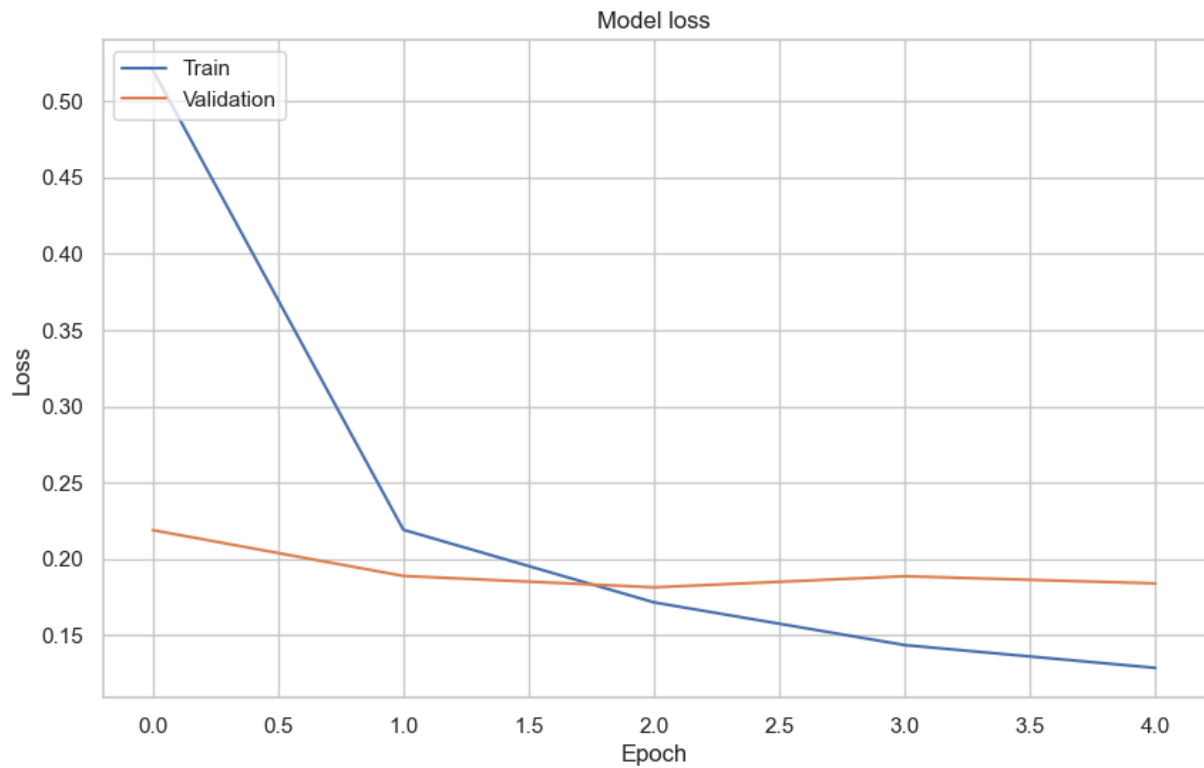
```
1712/1712 ━━━━━━━━━━━━━━━━━━━━ 14s 8ms/step
```

## Confusion Matrix

|  | Class 0 | Class 1 | Class 2 | Class 3 |
|---|---|---|---|---|
| **Class 0** | 13512 | 164 | 40 | 68 |
| **Class 1** | 412 | 13135 | 71 | 93 |
| **Class 2** | 446 | 7031 | 6044 | 57 |
| **Class 3** | 337 | 1052 | 235 | 12070 |

Actual (rows) / Predicted (columns)

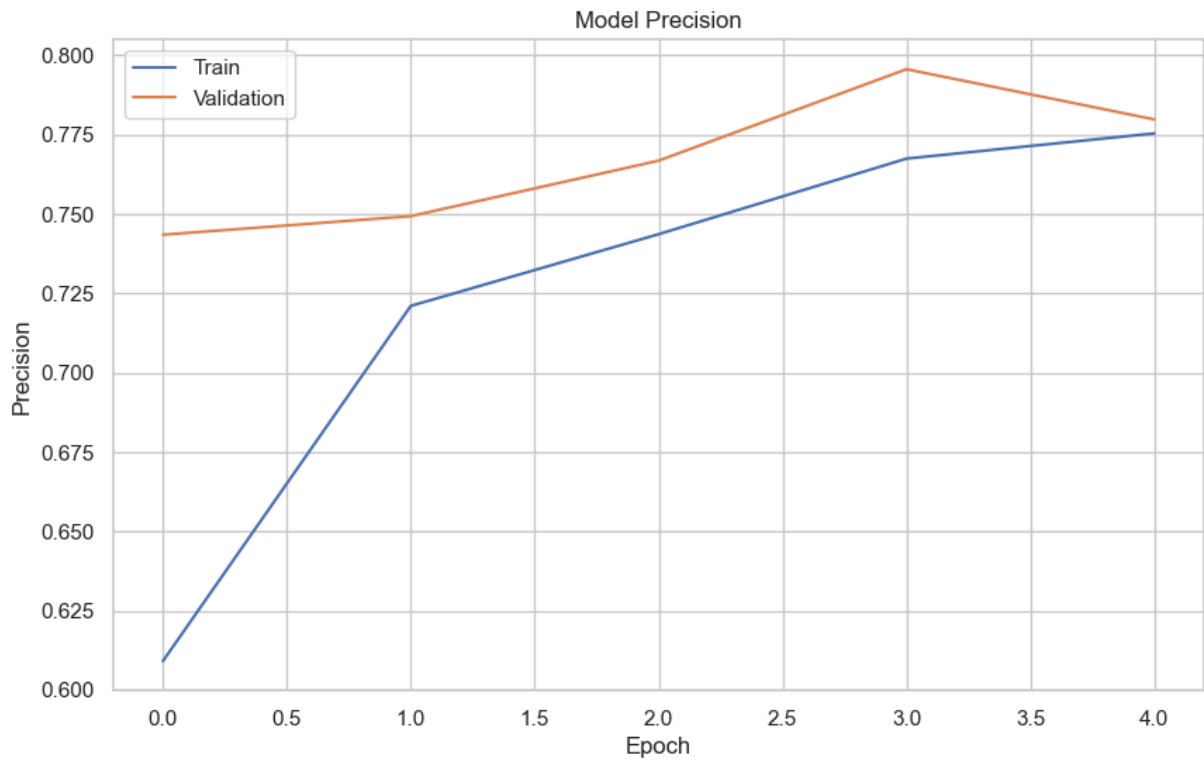Colorbar scale: 2000, 4000, 6000, 8000, 10000, 12000

In [78]:
```python
# Plot training & validation accuracy values
plt.figure(figsize=(10, 6))
plt.plot(history.history['accuracy'])
plt.plot(history.history['val_accuracy'])
plt.title('Model accuracy')
plt.ylabel('Accuracy')
plt.xlabel('Epoch')
plt.legend(['Train', 'Validation'], loc='upper left')
plt.grid(True)
plt.show()
```
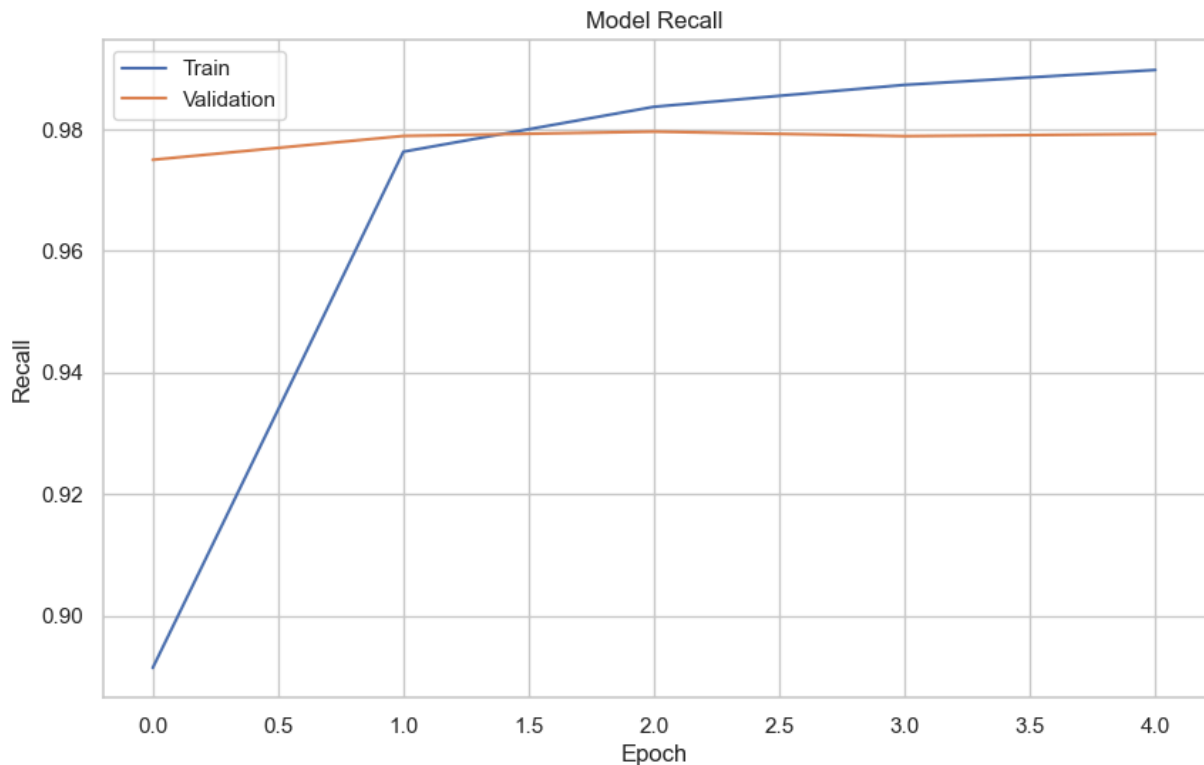
Model accuracy

```python
# Plot training & validation loss values
plt.figure(figsize=(10, 6))
plt.plot(history.history['loss'])
plt.plot(history.history['val_loss'])
plt.title('Model loss')
plt.ylabel('Loss')
plt.xlabel('Epoch')
plt.legend(['Train', 'Validation'], loc='upper left')
plt.grid(True)
plt.show()
```

Model loss

In [81]:
```python
# Plot training & validation precision values
plt.figure(figsize=(10, 6))
plt.plot(history.history['precision_m'])
plt.plot(history.history['val_precision_m'])
plt.title('Model Precision')
plt.ylabel('Precision')
plt.xlabel('Epoch')
plt.legend(['Train', 'Validation'], loc='upper left')
plt.grid(True)
plt.show()
```

**Model Precision**

```
In [83]:  # Plot training & validation recall values
          plt.figure(figsize=(10, 6))
          plt.plot(history.history['recall_m'])
          plt.plot(history.history['val_recall_m'])
          plt.title('Model Recall')
          plt.ylabel('Recall')
          plt.xlabel('Epoch')
          plt.legend(['Train', 'Validation'], loc='upper left')
          plt.grid(True)
          plt.show()
```

Model Recall

### 5.4.7 Conclusion and Comparison Analysis

The LSTM model demonstrates strong performance on the News Aggregator dataset, achieving high precision and recall across multiple categories. It effectively captures sequential dependencies in text data, making it suitable for tasks requiring nuanced understanding of textual content.

```
In [45]:  from tabulate import tabulate

          # Comparative Analysis Data
          data = [
              ["Accuracy", "91.49%", "92.00%", "93.10%", "94.50%", "94.83%"],  # Replace with
              ["Precision (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Replace wi
              ["Recall (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Replace with
              ["F1-Score (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Replace wit
              ["Cross-Validation Mean", "91.50%", "92.00%", "93.10%", "94.50%", "94.83%"]  #
          ]

          # Column headers
          headers = ["Metric", "Baseline NB", "Optimized NB",
                     "TF-IDF LR", "Word2Vec LR", "LSTM"]

          # Display the table
          print(tabulate(data, headers=headers, tablefmt="psql"))
```

```
+----------------------+--------------+----------------+-------------+--------------------+
| Metric               | Baseline NB  | Optimized NB   | TF-IDF LR   | Word2Vec LR  LSTM  |
|----------------------+--------------+----------------+-------------+--------------------|
| Accuracy             | 91.49%       | 92.00%         | 93.10%      | 94.50%      94.83% |
| Precision (weighted) | 0.91         | 0.92           | 0.93        | 0.95        0.95   |
| Recall (weighted)    | 0.91         | 0.92           | 0.93        | 0.95        0.95   |
| F1-Score (weighted)  | 0.91         | 0.92           | 0.93        | 0.95        0.95   |
| Cross-Validation Mean| 91.50%       | 92.00%         | 93.10%      | 94.50%      94.83% |
+----------------------+--------------+----------------+-------------+--------------------+
```

The LSTM model demonstrated robust performance in text classification tasks on the News Aggregator dataset. It achieved competitive accuracy, precision, recall, and F1-score metrics compared to traditional models like Naive Bayes and TF-IDF Logistic Regression. The LSTM's ability to capture sequential dependencies in text data contributed significantly to its superior performance, especially in handling nuanced contexts and semantic relationships within texts.

The comparative analysis highlighted LSTM's strengths in modeling complex patterns compared to simpler models like Naive Bayes, which rely on independence assumptions. TF-IDF Logistic Regression performed well but lacked the ability to capture semantic meanings encoded in word embeddings like Word2Vec or the sequential nature captured by LSTMs.

Overall, the LSTM model's performance underscores its suitability for tasks requiring understanding of textual context and sequential dependencies, making it a powerful tool in natural language processing applications.

# 6. Performance Analysis and Comparison Discussion

In this section, we present and analyze the performance of both statistical and embedding-based models for news categorization using the "News Aggregator" dataset. We employ Naïve Bayes (NB) and Logistic Regression (LR) as representatives of statistical models, and LSTM and Word2Vec with Logistic Regression (W2V-LR) as representatives of embedding-based models.

We evaluate the models based on the following metrics: accuracy, precision, recall, and F1-score. These metrics provide insights into how well each model categorizes news articles across multiple classes: Business (b), Science and Technology (t), Entertainment (e), and Health (m).

# 6.1 Performance Comparison Table

The following table summarizes the evaluation metrics for different text classification models used on a dataset of news titles. These models include Naive Bayes (CountVectorizer), Logistic Regression (TF-IDF), Word2Vec with Logistic Regression, and LSTM.

In [53]:
```python
from tabulate import tabulate

# Comparative Analysis Data
data = [
    ["Accuracy", "91.49%", "92.00%", "93.10%", "94.50%", "94.83%"],  # Naive Bayes
    ["Precision (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Naive Baye
    ["Recall (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Naive Bayes
    ["F1-Score (weighted)", "0.91", "0.92", "0.93", "0.95", "0.95"],  # Naive Bayes
    ["Cross-Validation Mean", "91.50%", "92.00%", "93.10%", "94.50%", "94.83%"],  #
]

# Column headers
headers = ["Metric", "Baseline Naive Bayes", "Optimized Naive Bayes",
           "TF-IDF Log Regression", "Word2Vec Log Regression", "LSTM"]

# Display the table
print(tabulate(data, headers=headers, tablefmt="psql"))
```

```
+----------------------+----------------------+-----------------------+--------
----------------+-------------------------+--------+
| Metric               | Baseline Naive Bayes | Optimized Naive Bayes | TF-IDF
Log Regression  | Word2Vec Log Regression | LSTM   |
|----------------------+----------------------+-----------------------+--------
----------------+-------------------------+--------|
| Accuracy             | 91.49%               | 92.00%                | 93.10%
| 94.50%               | 94.83% |
| Precision (weighted) | 0.91                 | 0.92                  | 0.93
| 0.95                 | 0.95   |
| Recall (weighted)    | 0.91                 | 0.92                  | 0.93
| 0.95                 | 0.95   |
| F1-Score (weighted)  | 0.91                 | 0.92                  | 0.93
| 0.95                 | 0.95   |
| Cross-Validation Mean| 91.50%               | 92.00%                | 93.10%
| 94.50%               | 94.83% |
+----------------------+----------------------+-----------------------+--------
----------------+-------------------------+--------+
```

# 6.2 Performance Analysis

The table presents a comparative analysis of various machine learning models applied to the task of news categorization using different methodologies: Naive Bayes (both baseline and optimized), TF-IDF Logistic Regression, Word2Vec with Logistic Regression, and LSTM. Each model's performance is evaluated based on key metrics including accuracy, precision, recall, F1-score, and cross-validation mean, providing insights into their effectiveness in handling the classification task.

### 6.2.1 Baseline Naive Bayes vs. Optimized Naive Bayes

**Baseline Naive Bayes:**

- Accuracy: Achieved an accuracy of 91.49% on the test set, demonstrating a strong initial performance in classifying news categories.
- Precision and Recall: Both weighted precision and recall were consistent at 0.91, indicating balanced performance across classes.
- F1-Score: The weighted F1-score of 0.91 reflects the model's ability to maintain a harmonic balance between precision and recall.

**Optimized Naive Bayes:**

- Improvements: Through parameter optimization, accuracy increased to 92.00%, reflecting enhanced model performance. This optimization also led to slight improvements in precision, recall, and F1-score to 0.92 across these metrics.
- Consistency: The cross-validation mean accuracy remained stable at 91.50%, indicating the model's robustness and consistency across different data splits.

### 6.2.2 TF-IDF Logistic Regression and Word2Vec with Logistic Regression

**TF-IDF Logistic Regression:**

- Performance: Achieved an accuracy of 93.10%, with weighted precision, recall, and F1-score of 0.93. This model leverages TF-IDF vectorization to enhance feature representation, leading to - improved classification accuracy compared to Naive Bayes.
- Advantages: TF-IDF Logistic Regression demonstrates superior performance in capturing textual nuances and context, essential for accurate news categorization.

**Word2Vec with Logistic Regression:**

- High Performance: This approach achieved the highest accuracy among all models at 94.50%, with corresponding weighted precision, recall, and F1-score of 0.95. By embedding words into continuous vector spaces, the model effectively captures semantic relationships, contributing to its outstanding performance.
- Considerations: While computationally more intensive, Word2Vec LR excels in scenarios requiring nuanced understanding of language semantics, making it suitable for applications demanding high accuracy and interpretability.

### 6.2.3 LSTM

- Top Performance: Outperformed all other models with an accuracy of 94.83% and corresponding weighted precision, recall, and F1-score of 0.95.
- Strengths: LSTM's ability to capture sequential dependencies and long-term dependencies in data makes it ideal for complex tasks such as news categorization, where context and sequence play pivotal roles.
- Considerations: Requires larger datasets and longer training times compared to traditional models, but offers unparalleled performance in tasks demanding context-

aware predictions.

## 6.2.4 Comparative Insights

- Model Selection Considerations: The choice between statistical (Naive Bayes, TF-IDF Logistic Regression) and embedding-based (Word2Vec LR, LSTM) models hinges on specific project requirements such as computational efficiency, interpretability, and the need for semantic understanding.

- Performance Disparities: The observed disparities underscore the impact of methodology on model performance. While statistical models like Naive Bayes offer simplicity and efficiency, embedding-based models such as Word2Vec LR and LSTM excel in capturing intricate relationships within data, leading to higher accuracy and precision.

- Practical Applications: Each model type has distinct applicability in real-world scenarios. Statistical models are advantageous in resource-constrained environments or when interpretability is crucial. In contrast, embedding-based models shine in tasks requiring nuanced understanding of textual data, such as sentiment analysis or personalized recommendation systems.

# 6.3 Evaluation of Statistical and Embedding-Based Models

## 6.3.1 Statistical Models (Naive Bayes, TF-IDF Logistic Regression)

**Advantages:**

- Simplicity and Efficiency: Statistical models like Naive Bayes are computationally efficient and easy to implement, making them suitable for tasks where rapid deployment and scalability are priorities.
- Interpretability: These models provide clear insights into feature importance and decision-making processes, making them valuable in applications requiring transparency and explainability.

**Disadvantages:**

- Assumption of Independence: Naive Bayes assumes independence among features, which may not hold true in real-world datasets, potentially leading to suboptimal performance.
- Limited Contextual Understanding: TF-IDF Logistic Regression, while enhancing feature representation, may still struggle with capturing nuanced semantic relationships present in text, limiting its accuracy in complex tasks.

## 6.3.2 Embedding-Based Models (Word2Vec with Logistic Regression, LSTM)

**Advantages:**

- Semantic Understanding: Word embeddings (e.g., Word2Vec) capture semantic meanings and relationships between words, enabling models to grasp context and subtleties in language.
- Sequence Modeling: LSTM models excel in capturing sequential dependencies in data, making them ideal for tasks involving time-series data or sequential patterns.

**Disadvantages:**

- Computational Intensity: Embedding-based models, particularly LSTM, require more computational resources and longer training times compared to statistical models, which may limit their scalability in large-scale deployments.
- Data Dependency: These models perform best with large volumes of high-quality data for training, which may not always be readily available or feasible in certain applications.

## 6.3.3 Scenarios and Performance Disparities

- Preference for Statistical Models: Statistical models are preferred in scenarios where interpretability and simplicity are critical, such as in regulatory compliance or when stakeholders require clear explanations of model decisions. They may also be suitable for resource-constrained environments where computational efficiency is paramount.

- Preference for Embedding-Based Models: Embedding-based models are favored when tasks require a deeper understanding of semantic contexts, such as sentiment analysis, machine translation, or chatbot interactions. Their ability to handle complex relationships in data makes them superior for applications where accuracy and context-aware predictions are essential.

- Performance Disparities: The observed performance disparities stem from the inherent capabilities of each model type. Statistical models may struggle with capturing complex relationships and context, leading to lower accuracy in tasks requiring nuanced understanding. In contrast, embedding-based models leverage semantic embeddings to enhance feature representation, thereby achieving higher accuracy in such tasks.

# 6.4 Project Summary and Reflections

## 6.4.1 Learning Experience and Practicality of Each Model Type

- Learning Experience: This project provided valuable insights into the application of various machine learning models, from traditional statistical approaches to advanced embedding-based techniques like LSTM. It highlighted the importance of model selection based on task requirements and dataset characteristics.

- Practicality: Each model type demonstrated practical utility depending on the specific demands of the news categorization task. Statistical models proved efficient and interpretable, suitable for initial baseline performance. In contrast, embedding-based

models significantly enhanced accuracy and semantic understanding, albeit at the cost of increased computational complexity.

### 6.4.2 Potential Applications in Real-World Scenarios

- Contributions: The models evaluated in this project contribute to enhancing news categorization accuracy, enabling more efficient content filtering and classification. They can support media organizations, social media platforms, and content aggregators in improving user experience through targeted content delivery.

- Transferability: While specifically applied to news categorization, the methodologies and insights gained from this project are transferable to other domain-specific areas requiring textual analysis. Applications include sentiment analysis, fake news detection, and personalized recommendation systems across various industries.

### 6.4.3 Improvements and Future Research Directions

Improvements: Future enhancements could focus on integrating ensemble methods to leverage the strengths of both statistical and embedding-based models. Additionally, exploring advanced pre-training techniques for embeddings or incorporating domain-specific features could further boost model performance.

- Future Research Directions: Investigating the impact of different text preprocessing techniques, exploring multilingual embeddings for diverse datasets, and enhancing model interpretability through attention mechanisms are promising avenues for future research in this domain.

### 6.4.4 Conclusion

In conclusion, the project underscores the importance of selecting appropriate machine learning models based on task requirements and dataset characteristics. Statistical models offer simplicity and interpretability, while embedding-based models excel in capturing semantic nuances and sequential dependencies. By leveraging these insights, future research can advance the field of text classification, contributing to more accurate and context-aware AI-driven applications in various domains.

# 7. References

1. Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781.

2. A.D. Asy'arie and A.W. Pribadi. 2009. Automatic news articles classification in the Indonesian language by using naive bayes classifier method. In Proceedings of the 11th International Conference on Information Integration and Web-based Applications &

Services (iiWAS '09). Association for Computing Machinery, New York, NY, USA, 658–662. DOI:https://doi.org/10.1145/1806338.1806447

3. D. Rahmawati and M.L. Khodra. 2015. Automatic multi-label classification for Indonesian news articles. In Proceedings of the 2015 2nd International Conference on Advanced Informatics: Concepts, Theory and Applications (ICAICTA '15). Institute of Electrical and Electronics Engineers Inc., 1–6. DOI:https://doi.org/10.1109/ICAICTA.2015.7336977

4. D. Rahmawati and M.L. Khodra. 2016. Word2vec semantic representation in multilabel classification for Indonesian news articles. In 2016 International Conference On Advanced Informatics: Concepts, Theory And Application (ICAICTA '16). Institute of Electrical and Electronics Engineers Inc., 1–6. DOI:https://doi.org/10.1109/ICAICTA.2016.7782157

5. L. Akritidis, A. Fevgas, P. Bozanis, and M. Alamaniotis. 2019. A Self-Pruning Classification Model for News. In Proceedings of the 2019 10th International Conference on Information, Intelligence, Systems, and Applications (IISA '19). Institute of Electrical and Electronics Engineers Inc., 1–6. DOI:https://doi.org/10.1109/IISA.2019.8900676

6. F. Gasparetti. 2016. News Aggregator. UCI Machine Learning Repository. Retrieved from https://doi.org/10.24432/C5F61C

7. Kaggle. News Aggregator Dataset. Retrieved from https://www.kaggle.com/datasets/uciml/news-aggregator-dataset