

# NLP Sentiment Analysis Assignment

## **Table of Contents**

<b>Abstract</b>	<b>1</b>
<b>1. Introduction and Motivation</b>	<b>1</b>
<b>2. Related Work</b>	<b>2</b>
<b>3. Experiment and Results</b>	<b>3</b>
<b>3.1 Preprocessing and Feature Selection:</b>	<b>3</b>
<b>3.1.1 Preprocessing</b>	<b>4</b>
<b>3.1.2 Feature Selection</b>	<b>7</b>
<b>3.2 Data Splits</b>	<b>9</b>
<b>3.3 Evaluation</b>	<b>10</b>
<b>3.4 Testing on Machine Learning Algorithms</b>	<b>10</b>
<b>3.4.1 Naïve Bayes</b>	<b>10</b>
<b>3.4.2 Logistic Regression and Support Vector Machines (SVMs)</b>	<b>15</b>
<b>3.5 Testing on Deep Learning Algorithm</b>	<b>19</b>
<b>3.5.1 BERT</b>	<b>19</b>
<b>4. Discussion</b>	<b>24</b>
<b>5. Conclusion and Future Work</b>	<b>25</b>
<b>References</b>	<b>25</b>

# Abstract

Sentiment Analysis is used for understanding the emotions of a given text and classifying them either into positive or negative reviews. The main aim of this project is to classify movie reviews from the IMDB dataset [1] as either positive or negative, by using four distinct algorithms: Naïve Bayes (implemented from scratch and using MultinomialNB), Support Vector Classifier (SVC), Logistics Regression, and BERT. For this classification task the steps taken are preprocessing and feature selection, including tokenization, N-grams, stopwords removal, lemmatization, stemming, and TF-IDF calculation. The models were evaluated using a classification report that included accuracy, precision, recall, and F1-score. Results show that the BERT model outperformed all the other algorithms, showing its effectiveness in sentiment analysis for movie reviews. The report gives a detailed explanation of each of the algorithms used and explains the reasons behind BERT's performance. In conclusion, this assignment not only identifies the best-performing model but also suggests some potential for future improvements.

## 1. Introduction and Motivation

Sentiment Analysis is the process of automatically detecting the sentiment in text (i.e., positive, negative, or neutral) [Spec].

It is an essential tool which is used by everyone, especially companies, to understand sentiments in all types of data. By automatically examining customer input, that is sentiments expressed in a survey, movie reviews or discussions, companies can gain insights of the text written and can understand what the reaction of the customer is. This information helps them to customize the product and services and helps them align with their customer's preferences and requirements [2].

Sentiment Analysis is used by many top companies like Apple, Google, KFC and many more. One of the use cases of Google using sentiment analysis is to make Chrome better. They check on what users want to say, not just words, but also check if it is positive or negative and what they want to say about different parts of Chrome. This helps them see what is good or bad and helps them make Chrome better. [3]

Even though it a great tool for automatically detecting sentiment in a text there are limitations, some of them are [4]: Context sensitive error – trying to express a negative statement but using positive wordings (e.g.: "Sure, what a brilliant suggestion!"). This phrase can be treated as positive or negative and sentiment analysis might wrongly classify them as positive instead of negative depending on the situation. Multilingual data – although English is the language used worldwide since sentiment analysis is used by many companies worldwide the reviews or feedback received from customers maybe of a different language. Since sentiment analysis is mainly trained in English, some of the data might get lost or mismatched during translations and it will not provide the correct analysis. Emojis – In this century everyone responds or replies by using emojis. Sentiment Analysis is a text-based tool which depends on the text written to classify it into positive or negative. If the text consists of emojis then that text will not be analysed [4]. These are some of the limitations or challenges faced by sentiment analysis.

Different Machine Learning and Deep Learning techniques are used for Sentiment Analysis. Machine Learning techniques uses classification algorithms like Naïve Bayes, Support Vector

Machines (SVMs), Logistic Regression, Linear Regression to do Sentiment Analysis [5]. Deep Learning techniques transform the text into embeddings and some of the algorithms include Long Short-Term Memory (LSTM) and other recurrent neural networks (RNNs), Convolutional Neural Networks (CNNs) and Transformers (BERT and GPT-3) [6,7].

In this project I am testing to see which algorithm is better for sentiment analysis. The IMDB Reviews data [1] have been used for sentiment analysis and have been tested on 4 different algorithms (Naïve Bayes (from scratch and MultinomialNB), Logistic Regression, SVM, and BERT). Different feature selection and feature generation methods were applied on the dataset and different combinations of features were tested on these algorithms.

Different strategies were used and tested to see which works best. For example, Logistic Regression and SVM was first tested using default hyperparameters and then hyperparameter optimizations was done similarly in the case of Naïve Bayes it was tested using the built-in version and scratch version. Out of all the models tested the BERT model performed better than other models.

## 2. Related Work

There are a several existing methods for Sentiment Analysis some of them include [5,6,7]:

- **Machine Learning methods** like Naïve Bayes [10], Support Vector Machines (SVMs) [12], Random Forests (RFs) [13] and many more.
  - Naïve Bayes (NB) – In this paper [9], Naïve Bayes algorithm was used to do sentiment analysis on Twitter tweets. That is classifying tweets into positive and negative. Multinomial NB is a probabilistic learning method. The goal of this method is to find the best class in a document (Maximum Posterior class (MAP) or most likely class) [10].
  - Support Vector Machines (SVMs) – It is a binary classification algorithm. It works by transforming data into a high-dimensional space and drawing a line to separate them into two groups. It works well when there is a limited amount of data provided to the model [11,12,18].
  - Random Forests (RF) – It is a machine learning method which works based on the principle of several decision trees put together into one to make predictions. It works for both classification and regression tasks. The larger the number of decisions trees put together, the better will be the performance given by Random Forest [13,14].
- **Deep Learning methods** include Long Short-Term Memory (LSTM) [15], Recurrent Neural Networks (RNNs) [16], and Transformer based models like BERT [17] and many more.
  - Long Short-Term Memory (LSTM) – It is a type of neural network used in Deep Learning. It is great at learning from sequences of data, like speech recognition and translation. It is a special kind of Recurrent Neural Network (RNN) that works well for different tasks. There are two problems with RNNs vanishing gradient and exploding gradient, but these two problems are solved by LSTMs [15,19]
  - Recurrent Neural Networks (RNNs) – It is a type of artificial neural network which uses sequential data. Unlike other methods, RNN remember what they have seen before and use that to understand what is happening now. They cannot predict the future

because they only look at the past. It is used in translation, speech recognition and voice assistants [16, 20].

- BERT (Bidirectional Encoder Representations from Transformers) – “It is transformers based Deep Learning Model where each output element is linked to the input element and the weights are calculated dynamically between them.” [21,17].
- VADER – A Rule-based model for Sentiment Analysis of social media Text. This model was compared with eleven other methods (i.e., LIWC, ANEW, the General Inquirer, SentiWordNet, and machine learning oriented techniques relying on Naive Bayes, Maximum Entropy, and Support Vector Machine (SVM) algorithms) and it was found that VADER works better than most of these models and performed better than humans who had classified these texts manually. Two strategies were used in this method, first they made a list of words on how people feel in the online text messages and figured out the strength of these words in showing those feelings. Then made five simple rules on how these words are used to show feelings [8].

### 3. Experiment and Results

In the CodeA folder there are three .ipynb files, three .csv files for BERT and data folder with pos and neg reviews:

1. Codes for experiment and results of Naïve Bayes (scratch), MultinomialNB, SVCs and Logistic Regression can be found in **CodeA – Part1.ipynb**
2. Code for experimenting different features can be found in **Feature Test.ipynb**
3. Code for BERT model can be found in **BERT.ipynb**

#### 3.1 Preprocessing and Feature Selection:

These are the list of libraries which have been imported for running each of the techniques mentioned below:

```
1 #Libraries
2
3 import numpy as np
4 import pandas as pd
5 import nltk
6 from nltk import ngrams
7
8 import string
9 import math
10
11 nltk.download('stopwords')
12 nltk.download('wordnet')
13
14 stoplist = set(nltk.corpus.stopwords.words('english'))
15
16 from nltk.corpus import stopwords
17 from nltk.tokenize import word_tokenize
18 from nltk.stem.lancaster import LancasterStemmer
19 from sklearn.feature_extraction.text import TfidfVectorizer
20 from nltk.stem import WordNetLemmatizer
21 from sklearn.model_selection import train_test_split
22 from sklearn.naive_bayes import MultinomialNB
23 from sklearn.svm import SVC
24 from sklearn.linear_model import LogisticRegression
```

*Figure 1: List of Libraries*

Preprocessing and Feature Selection of the data are the important steps for carrying out any NLP tasks. The different ways are [22]:

- Tokenization
- N-Grams
- Stopwords
- Lemmatization
- Stemming

### 3.1.1 Preprocessing

1. **Tokenization**<sup>1</sup> - is the process of dividing continuous textual information into words, phrases, sentences, symbols, and other components known as tokens [22]. In the data provided, I have converted the whole 'Reviews' into tokens.

```
1 #Tokenization - Data
2 from nltk.tokenize import word_tokenize
3
4 tokenized = []
5
6 for i in df.index:
7     tokenized.append(word_tokenize(df['reviews'][i].lower()))
8     #tokenized.sort()
9
10 df.insert(3, 'tokenized_review', tokenized, True)
11 df.head()
```

Figure 2: Tokenization of Reviews

The tool used for tokenization is - from nltk.tokenize import word\_tokenizer

2. **Frequency Distribution**<sup>2</sup> – This was used to check the frequency of words in a particular document. I have tested to see which words are frequent in two of the reviews text files and plotted a graph to see the frequencies. From the graph it is evident that Stopwords and Punctuation are more frequent in two of the text files. This means even though they are more frequent they do not carry any useful meanings. So, they have been removed in the Feature Selection part.

```
frequency_words = nltk.FreqDist(tokenized[0])
print(frequency_words.most_common(50))
```

---

<sup>1</sup> [https://github.com/KuWathsala/NLP-Sentiment-Analysis-Beginners-/blob/master/NLP\\_SA.ipynb](https://github.com/KuWathsala/NLP-Sentiment-Analysis-Beginners-/blob/master/NLP_SA.ipynb)  
[https://colab.research.google.com/drive/1C5RLXOvaKgbdmeW4PPjg\\_Q4xooX-r4Oz](https://colab.research.google.com/drive/1C5RLXOvaKgbdmeW4PPjg_Q4xooX-r4Oz)  
<https://colab.research.google.com/drive/1WJ2HEQtwpq7TsGRyV5E1tO6IFGnCbW9A?usp=sharing>

<sup>2</sup> [https://colab.research.google.com/drive/1mkYrOnFt60iDzcH0bhOIP\\_Us0dvDmlt?usp=sharing#scrollTo=TjTSsQxCvNIF](https://colab.research.google.com/drive/1mkYrOnFt60iDzcH0bhOIP_Us0dvDmlt?usp=sharing#scrollTo=TjTSsQxCvNIF)

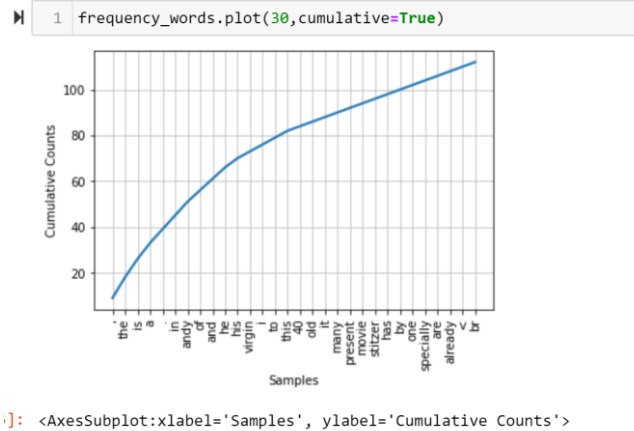


Figure 3: Checking Frequency Distribution of Tokenized first document.

The tool used for Frequency Distribution – `nltk.FreqDist()`

3. **N-Grams**<sup>3</sup> – It refers to the series on N-words. There are different types of N-Grams that are Unigram, Bigrams, Trigrams and more. In N-1 word contexts, N-gram models predict the most frequently appearing words that come after the given sentence. This model is a probabilistic model trained on several texts [23]. For example, sentence - 'This is an assignment.'

1. Unigram = ['This', 'is', 'an', 'assignment']
2. Bigrams = ['This is', 'is an', 'an assignment']
3. Trigrams = ['This is an', 'is an assignment', '... ...']
4. Similarly, this applies for other grams as well.

```

from nltk import ngrams
terms = list(ngrams(doc[0].split(), n))

```

Figure 4: Applying N-grams on the reviews (features /terms)

The tool used for N-Grams is - `from nltk import ngrams`.

On this movie review dataset, I have tested using Unigram, Bigrams and Trigrams. All the grams were applied on the tokenized reviews.

4. **Stopwords**<sup>4</sup> – It is a collection of frequently used words in the language. That is, 'a', 'is', 'an', 'the' and many more. These words are often removed because they usually carry a limited amount of information [24].

The tool used for Stopwords is - `from nltk.corpus import stopwords`.

<sup>3</sup> <https://spotintelligence.com/2023/04/05/n-grams/>

<sup>4</sup> <https://colab.research.google.com/drive/1IsqgtoEVG8n21tb0tUMCN0mLqiBbgnJ8> and from previous Coursework done in Machine Learning.

5. **Stemming**<sup>5</sup> – It is a text normalization method used in NLP, which converts the text to its base form. For example, 'rain' will be stem for 'raining', 'rained', 'rains' [25]. It removes the final characters of a word, which can sometimes result in inaccuracies in meanings and spellings (i.e., for 'caring' the stemmed word will be 'car'). Stemming is usually used when there is large amount of data where performance is a concern [26].

```
# Apply stemming
st = LancasterStemmer()
stemmed_documents_1 = [[' '.join([st.stem(word) for word in words])] for words in documents]
```

Figure 5: Applying Stemming to the words in the documents.

The tool used for Stemming is - `from nltk.stem.lancaster import LancasterStemmer`.

6. **Lemmatization**<sup>6</sup> - It is also a text normalization method used in NLP and works almost like Stemming. It considers the context and transforms the word into its base form known as 'lemma'. (i.e., for 'caring' the lemma would be correctly identified as 'care'). It is computationally very expensive since it has lookup tables and other features [26].

```
# Apply Lemmatization
lemmatizer = WordNetLemmatizer()
lemmatized_documents_1 = [[' '.join([lemmatizer.lemmatize(word) for word in words if word not in
                                     stoplist and word not in string.punctuation])] for words in documents]
```

Figure 6: Applying Lemmatization to the words in the documents

The tool used for Lemmatization is - `from nltk.stem import WordNetLemmatizer`.

7. **TF-IDF**<sup>7</sup> – A statistical method which is used to evaluate the relevancy of a word is to a document in the collections of documents. The TF\_IDF of a word is calculated by multiplying the Term Frequency of word in the document and Inverse Document Frequency of word across a set of documents [27].

### **Why TF-IDF was chosen over other?**

There are three different kinds of Normalization techniques: TF-IDF, Term Frequency and PPMI (Positive Pointwise Mutual Information).

Out of these three techniques the one I chose for normalizing the movie reviews was TF-IDF. This is because TF-IDF normalizes term frequencies based on the importance of not just current document but across the entire collection of documents in this case of movie reviews. TF-IDF gives higher weights to terms that are rare across the entire collection of documents, making it more effective in identifying important terms [29,30]. PPMI takes semantic relations between terms. It is used in distributional semantics to measure association between words based on their occurrence in the context of document. PPMI is used to capture semantic relationships between terms rather than sentiment [31]. So, the widely used approach for sentiment analysis is TF-IDF as it

---

<sup>5</sup> <https://colab.research.google.com/drive/1sqqt0EVG8n21tb0tUMCN0mLqiBbgnJ8> and from previous Coursework done in Machine Learning from Ekaterina.

<sup>6</sup> <https://www.geeksforgeeks.org/python-lemmatization-with-nltk/>

<sup>7</sup> <https://colab.research.google.com/drive/1sqqt0EVG8n21tb0tUMCN0mLqiBbgnJ8> and from previous coursework done in Machine Learning from Ekaterina



takes importance of words within a text by considering the entire document and helps in identifying if the phrase is in positive or negative sentiment [32].

### 3.1.2 Feature Selection

Within the Feature Test.ipynb several features have been tested on Unigrams without threshold, Bigrams and Trigrams with Threshold and are applied to both Stemmed and Lemmatized words. Of the ones tested the three features below are the ones which gave a higher accuracy when tested using Multinomial NB. So, these were taken as three main features for testing other models.

Three sets of Features have been developed with a combination of Stopwords, Stemming, Lemmatization and N-grams.

#### Feature 1 – Stemming, Bigrams with threshold and Applied TF-IDF

```
1 # Function to calculate TF-IDF with bi-grams and
2 # filter based on frequency threshold and after applying stemming
3 def calculate_tfidf_ngrams_with_threshold(documents, n, threshold):
4     # Apply stemming
5     st = LancasterStemmer()
6     stemmed_documents_1 = [[' '.join([st.stem(word) for word in words])] for words in documents]
7
8     # Generate n-grams and calculate frequency
9     all_bigrams = []
10    for doc in stemmed_documents_1:
11        # Splitting terms in the document to words (bigrams)
12        terms = list(ngrams(doc[0].split(), n))
13        all_bigrams.extend(terms)
14
15    # Calculate the frequency of all bigrams
16    frequency_word_bigrams = nltk.FreqDist(all_bigrams)
17
18    # Filter bigrams based on the threshold
19    filtered_bigrams = {bigram: freq for bigram, freq in frequency_word_bigrams.items() if freq >= threshold}
20
21    # Create a vocabulary
22    all_terms = set(filtered_bigrams.keys())
23    vocabulary_1 = sorted(list(all_terms))
24
25    # Calculate TF - Asked ChatGPT on how to create a TF-matrix so it suggested using pd.DataFrame
26    # https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.html
27    tf_matrix = pd.DataFrame(0, index=df.index, columns=vocabulary_1, dtype=float)
28    # Iterating through index and documents
29    for i, doc in enumerate(stemmed_documents_1): # https://realpython.com/python-enumerate/
30        terms = list(ngrams(doc[0].split(), n))
31        for term in terms:
32            # counting the frequency of each term in each document.
33            # Incrementing the tf matrix by 1 if the term is in the vocabulary
34            if term in vocabulary_1:
35                tf_matrix.at[i, term] += 1
36
37    # Calculate IDF - Asked ChatGPT on how to create a IDF-vectors so it suggested using pd.Series
38    # https://pandas.pydata.org/docs/reference/api/pandas.Series.html
39    idf_vector = pd.Series(0, index=vocabulary_1, dtype=float)
40    N = len(stemmed_documents_1)
41    for term in vocabulary_1:
42        # counting the number of documents in which the bigram appears
43        df_term = sum([1 for doc in stemmed_documents_1 if term in list(ngrams(doc[0].split(), n))])
44        # Calculating the IDF value of each bigram
45        idf_vector.at[term] = math.log(N / (1 + df_term), 10)
46
47    # Calculate TF-IDF
48    tfidf_matrix_1 = tf_matrix * idf_vector
49
50    return tfidf_matrix_1, vocabulary_1
51
52 # Calculate TF-IDF matrix with bigrams and filter based on frequency threshold
53 tfidf_matrix_1, vocabulary_1 = calculate_tfidf_ngrams_with_threshold(df['tokenized_review'], n=2, threshold=100)
54
```

Figure 7: Implementation of Feature 1

In this all the tokenized words of the documents are first stemmed. Bigrams are applied to the terms in those stemmed documents. After that filtering of the bigram is done so that bigrams which are less frequent are removed. A threshold of 100 is set, this means any bigram which has frequency less than 100 is not included in the vocabulary list. TF-IDF calculation is then done to this set of features. TF matrix is created to calculate the TF and IDF vector to calculate the IDF. Then both are multiplied to get the TF-IDF matrix which is used as transformed feature data for the Data split.

TF-Matrix is calculated by taking all the documents in the stemmed\_document and again applying Bigrams to that and then checking to see if those terms are there in the vocabulary (filtered\_bigrams) and incrementing the matrix by 1 if the term exists. That is counting the frequency of each term (bigram) in the document.

IDF-Vector is calculated by taking the term from the vocabulary and then counting the number of documents in which this bigram appears. Then it calculates the IDF value of that current bigram using the formula  $(\log(N/1 + df\_term), 10)$  where N is total number of documents and df.term is document frequency of the bigram.

Then both the TF-Matrix and IDF-Vector is then multiplied to get the TF-IDF matrix

Code for calculation of TF-IDF is taken from the Practice Labs provided and changes have been made accordingly.

## Feature 2 – Stemming, Unigrams without threshold, removing Stopwords and punctuations and Applied TF-IDF

```

1 # Function to calculate TF-IDF after applying stemming
2 # and stopwords and punctuations removal on unigrams
3 def calculate_tfidf(documents):
4     # Apply stemming
5     st = LancasterStemmer()
6     stemmed_documents_2 = [[' '.join([st.stem(word) for word in words if word not in
7                                     stoplist and word not in string.punctuation]]) for words in documents]]
8
9
10    # Create a vocabulary
11    all_terms = set() #only unique terms are allowed
12    for doc in stemmed_documents_2:
13        all_terms.update(doc[0].split()) #.update() to add items to dictionary
14    vocabulary_2 = sorted(list(all_terms))
15
16    # Calculate TF - Asked ChatGPT on how to create a TF-matrix so it suggested using pd.DataFrame
17    tf_matrix = pd.DataFrame(0, index=df.index, columns=vocabulary_2, dtype=float)
18    #Iterating through index and document
19    for i, doc in enumerate(stemmed_documents_2): #https://realpython.com/python-enumerate/
20        #counting the frequency of each term in each document.
21        #Incrementing the tf matrix by 1 if the term is in the document
22        for term in doc[0].split():
23            tf_matrix.at[i, term] += 1 #adding to the matrix
24
25    # Calculate IDF - Asked ChatGPT on how to create a IDF-vectors so it suggested using pd.Series
26    # https://pandas.pydata.org/docs/reference/api/pandas.Series.html
27    idf_vector = pd.Series(0, index=vocabulary_2, dtype=float)
28    N = len(stemmed_documents_2)
29    for term in vocabulary_2:
30        #counting the number of documents in which the term appears
31        df_term = sum([1 for doc in stemmed_documents_2 if term in doc[0].split()])
32        #Calculating the IDF value of each unigram
33        idf_vector.at[term] = math.log(N / (1 + df_term), 10)
34
35    # Calculate TF-IDF
36    tfidf_matrix_2 = tf_matrix * idf_vector
37
38    return tfidf_matrix_2, vocabulary_2
39
40    # Calculate TF-IDF matrix and get vocabulary
41    tfidf_matrix_2, vocabulary_2 = calculate_tfidf(df['tokenized_review'])
42
43

```

Figure 8: Implementation of Feature 2

A similar approach has been taken for generating this feature, the only differences are it is not considering the threshold, so this means that the vocabulary is created by splitting all the stemmed words of the document (unigram) and the Stopwords and Punctuations are removed. And then again TF-IDF matrix is created as explained above without any n-grams splitting but instead all the words in the document are split and compared against the terms in each document. And then IDF is calculated to count the number of documents in which that term appears. One of the drawbacks of this method is that since it's not considering any threshold it will be checking against all the terms within documents which is time consuming. So, as a future development I will incorporate threshold into this method.

### Feature 3 – Lemmatization, Unigrams without threshold, removing Stopwords and punctuations and Applied TF-IDF

```
1 # Function to calculate TF-IDF after applying Lemmatization
2 # and stopwords and punctuations removal on unigrams
3 def calculate_tfidf(documents):
4     # Apply Lemmatization
5     lemmatizer = WordNetLemmatizer()
6     lemmatized_documents_1 = [[' '.join([lemmatizer.lemmatize(word) for word in words if word not in
7                                         stoplist and word not in string.punctuation])] for words in documents]
8
9
10    # Create a vocabulary
11    all_terms = set() # Only unique terms are allowed
12    for doc in lemmatized_documents_1:
13        all_terms.update(doc[0].split()) #.update() to add items to dictionary
14    vocabulary_3 = sorted(list(all_terms))
15
16    # Calculate TF - Asked ChatGPT on how to create a TF-matrix so it suggested using pd.DataFrame
17    tf_matrix = pd.DataFrame(0, index=df.index, columns=vocabulary_3, dtype=float)
18    # Iterating through index and document
19    for i, doc in enumerate(lemmatized_documents_1): #https://realpython.com/python-enumerate/
20        # counting the frequency of each term in each document.
21        # Incrementing the tf matrix by 1 if the term is in the document
22        for term in doc[0].split():
23            tf_matrix.at[i, term] += 1 #adding to matrix
24
25    # Calculate IDF - Asked ChatGPT on how to create a IDF-vectors so it suggested using pd.Series
26    # https://pandas.pydata.org/docs/reference/api/pandas.Series.html
27    idf_vector = pd.Series(0, index=vocabulary_3, dtype=float)
28    N = len(lemmatized_documents_1)
29    for term in vocabulary_3:
30        # counting the number of documents in which the term appears
31        df_term = sum([1 for doc in lemmatized_documents_1 if term in doc[0].split()])
32        # Calculating the IDF value of each unigram
33        idf_vector.at[term] = math.log(N / (1 + df_term), 10)
34
35    # Calculate TF-IDF
36    tfidf_matrix_3 = tf_matrix * idf_vector
37
38    return tfidf_matrix_3, vocabulary_3
39
40 # Calculate TF-IDF matrix and get vocabulary
41 tfidf_matrix_3, vocabulary_3 = calculate_tfidf(df['tokenized_review'])
42
```

Figure 9: Implementation of Feature 3

This approach is very similar to the Feature 2 development, the only difference is that it does Lemmatization instead of Stemming.

## 3.2 Data Splits<sup>8</sup>

Sklearn's `train_test_split` [28] function was used to split the features extracted into train, development (dev), and test sets.

First the vectorized feature is divided into train and test sets with train set having 80% and test set having 20% of the whole data.

Then, the train set is further divided into train and development set with train set having 80% and development set having 20% of the further divided train data.

```
Split data into train, development and test sets with Feature 1 - Stemming, bigrams and Applied TF-IDF
1 X_train1, X_test1, y_train1, y_test1 = train_test_split(tfidf_matrix_1, df.labels, test_size=0.2, random_state=42)
2 X_train1, X_dev1, y_train1, y_dev1 = train_test_split(X_train1, y_train1, test_size=0.2, train_size=0.8, random_state=42)

Split data into train, development and test sets with Feature 2 - Stemming, removed stopwords and punctuations, unigrams and TF-IDF
1 X_train2, X_test2, y_train2, y_test2 = train_test_split(tfidf_matrix_2, df.labels, test_size=0.2, random_state=42)
2 X_train2, X_dev2, y_train2, y_dev2 = train_test_split(X_train2, y_train2, test_size=0.2, train_size=0.8, random_state=42)

Split data into train, development and test sets with Feature 3 - Lemmatization, removed stopwords and punctuations, unigrams and TF-IDF
1 X_train3, X_test3, y_train3, y_test3 = train_test_split(tfidf_matrix_3, df.labels, test_size=0.2, random_state=42)
2 X_train3, X_dev3, y_train3, y_dev3 = train_test_split(X_train3, y_train3, test_size=0.2, train_size=0.8, random_state=42)
```

Figure 10: Splitting of data

The tool used for Data Splits is - from `sklearn.model_selection` import `train_test_split`

<sup>8</sup> <https://realpython.com/train-test-split-python-data/>

### 3.3 Evaluation<sup>9</sup>

Evaluation on the Test set for all the models was done by checking the Accuracy, Precision, Recall and F1 Score [40]:

- 1) Classification report – It shows the Precision, Recall, F1-Score, Support and Accuracy that is related to the model tested.
- 2) Precision – Calculates how exactly the classifier has made the decisions.
- 3) Recall – Measure the ability of the classifier to correctly identify all the positive values.
- 4) F1- Score – It calculates the harmonic mean of the Precision and Recall.
- 5) Support – Number of occurrences of that class in the dataset.
- 6) Accuracy – Shows the percentage of how well the model has performed.
- 7) Confusion Matrix – Measures the performance of the classification model. Consists of TP (True Positives), FP (False Positives), TN (True Negatives) and FN (False Negatives). TP – Predicts all the positive reviews correctly, TN – Predicts all negative reviews correctly, FP – Predicted some negative reviews as positive, FN – Predicted some positive reviews as negative [41].

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figure 11: Confusion Matrix - (Image courtesy: My Photoshopped Collection) [41]

### 3.4 Testing on Machine Learning Algorithms

#### 3.4.1 Naïve Bayes

**Naïve Bayes** is a classification algorithm which uses the probabilistic approach. The basis of Naïve Bayes comes from the Bayes Theorem [33].

According to Bayes Theorem [33]:

$$P(C|x) = \frac{P(x|C) P(C)}{P(x)}$$

That is:

P (C | x) - probability of finding a class given the observation (Posterior Probability)

P (x | C) - probability of finding the observation given a particular class (Likelihood)

P (C) – probability of class (Prior Probability)

P(x) – probability of observations

---

<sup>9</sup> <https://github.com/ReiCHU31/Sentiment-analysis-of-IMDb-movie-reviews/blob/master/Sentiment%20analysis%20of%20IMDb%20movie%20reviews.ipynb>

### 3.4.1.1 Naïve Bayes from Scratch<sup>10</sup>

Naïve Bayes algorithm calculates the posterior probability of each class given the observation and assigns the prediction to the class with the highest posterior probability [33].

The code for the implementation of the Naïve Bayes from scratch was **adapted from** [33] and with the help of ChatGPT.

```
1 from sklearn.metrics import accuracy_score
2 from sklearn.preprocessing import LabelEncoder
3
4 # Train the Naive Bayes model
5 def train_naive_bayes(X_train, y_train):
6     # Calculate class priors
7     prior = {}
8     classes, counts = np.unique(y_train, return_counts=True)
9     # calculate P(class) = count(class 0 or 1)/total number of values in the class
10    # Prior probability
11    for class_0_1, count in zip(classes, counts):
12        prior[class_0_1] = count / len(y_train)
13
14    # Calculate means and variances for each feature in each class
15    means = {}
16    var = {}
17
18    # Taking feature of each of the class - (x)
19    for cls in classes:
20        class_indices = [i for i in range(len(y_train_encoded)) if y_train_encoded[i] == cls]
21        #https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.iloc.html
22        class_data = X_train.iloc[class_indices] #Access the feature of a particular class using iloc
23
24        class_means = np.mean(class_data, axis=0)
25        class_var = np.var(class_data, axis=0)
26
27        means[cls] = class_means
28        var[cls] = class_var
29
30    return prior, means, var
31
32
33 #Gaussian Distribution
34 def Normal(n, mu, var):
35     sd = np.sqrt(var)
36     # Check if standard deviation is zero
37     if sd == 0:
38         return np.zeros_like(n) # Return an array of zeros to avoid division by zero
39
40     pdf = (np.e ** (-0.5 * ((n - mu)/sd) ** 2)) / (sd * np.sqrt(2 * np.pi))
41     return pdf
42
43 # Prediction of Naive Bayes
44 def predict_naive_bayes(X, prior, means, var):
45     classes = list(prior.keys())
46     predictions = []
47
48     # calculate P(reviews|class)
49     for i in range(len(X)):
50         class_likelihoods = []
51
52         for cls in classes:
53             feature_likelihoods = []
54             # Log-prior - P(class)
55             feature_likelihoods.append(np.log(prior[cls]))
56
57             for j in range(X.shape[1]): # Use shape[1] to iterate over features
58                 data = X.iloc[i, j] # Access the j-th feature of the i-th sample using iloc
59                 mean = means[cls][j]
60                 variance = var[cls][j]
61                 #Likelihoods - p(reviews1|class) + p(reviews2|class) + ... + p(reviewsN|class)/P(class)
62                 likelihood = Normal(data, mean, variance)
63                 # Log-Likelihood
64                 if likelihood != 0:
65                     likelihood = np.log(likelihood)
66                 else:
67                     likelihood = 1 / X.shape[1] # Adjust the handling of zero likelihood
68
69                 feature_likelihoods.append(likelihood)
70
71             total_likelihood = sum(feature_likelihoods)
72             # Calculate posterior - P(class|reviews)
73             class_likelihoods.append(total_likelihood)
74
75         # Take class with max log-likelihood
76         max_index = class_likelihoods.index(max(class_likelihoods))
77         prediction = classes[max_index]
78         predictions.append(prediction)
79
80     return predictions
81
```

Figure 12: Implementation of Naive Bayes from scratch

<sup>10</sup> <https://blog.devgenius.io/implementing-na%C3%AFve-bayes-classification-from-scratch-with-python-badd5a9be9c3>

There are three functions in the code (`train_naive_bayes ()`, `Normal ()` and `predict_naive_bayes ()`). The `train_naive_bayes ()` function calculates the prior probability ( $P(\text{class})$ ) of the class and then calculates the mean and variance of the features of a particular class.

`Predict_naive_bayes` function is used to predict the posterior probability when the dev and test datasets are given.

As explained before to calculate the posterior probability ( $P(\text{class} | \text{reviews})$ ) we need calculate the likelihood ( $P(\text{reviews} | \text{class})$ ) and prior  $P(\text{class})$  probabilities [33]. To calculate the likelihood of the distribution of reviews the method used here is Gaussian distribution. The function `Normal ()` calculates the gaussian distribution given the mean and variance of the features of a particular class.

To calculate the log-likelihood, steps are taken which calculates the `np.log(likelihood)` if the likelihood calculated before using Gaussian distribution is not 0 and `1/X.Shape [1]` (1/number of columns (features) in the array) if the likelihood is 0.

All the `np.log (prior [class])` and `np.log(likelihood)` values are then added to the `feature_likelihood` list. And they are all summed together and added into the `class_likelihood` list. From that `max_index` takes the class with maximum `likelihood` and prediction of the class is done in that way to decide if it a positive or negative review (feature).

A Label Encoder was used to convert all the Labels (classes) into numerical format so it can be used in the Naïve Bayes classifier above.

```
1 label_encoder = LabelEncoder()

# Encode class labels
y_train_encoded = label_encoder.fit_transform(y_train1)
y_dev_encoded = label_encoder.transform(y_dev1)
y_test_encoded = label_encoder.transform(y_test1)
```

Figure 13: Label Encoder on the Class Labels

Evaluation of this model was done on all the three features created above:

Feature 1

```
1 # Encode class labels
2 y_train_encoded = label_encoder.fit_transform(y_train1)
3 y_dev_encoded = label_encoder.transform(y_dev1)
4 y_test_encoded = label_encoder.transform(y_test1)
5
6 #Calculate prior, mean and variance of trainset
7 prior, means, var = train_naive_bayes(X_train1, y_train_encoded)
8
9 # Make predictions on the development set
10 dev_predictions = predict_naive_bayes(X_dev1, prior, means, var)
11
12 # Evaluate the model on the development set
13 dev_accuracy = accuracy_score(y_dev_encoded, dev_predictions)
14 print(f"Development Set Accuracy1: {dev_accuracy}")

Development Set Accuracy1: 0.728125
```

#### Feature 2

```
1 # Encode class Labels
2 y_train_encoded = label_encoder.fit_transform(y_train2)
3 y_dev_encoded = label_encoder.transform(y_dev2)
4 y_test_encoded = label_encoder.transform(y_test2)
5
6 #Calculate prior, mean and variance of trainset
7 prior, means, var = train_naive_bayes(X_train2, y_train_encoded)
8
9 # Make predictions on the development set
10 dev_predictions = predict_naive_bayes(X_dev2, prior, means, var)
11
12 # Evaluate the model on the development set
13 dev_accuracy = accuracy_score(y_dev_encoded, dev_predictions)
14 print(f'Development Set Accuracy2: {dev_accuracy}')
```

Development Set Accuracy2: 0.55

#### Feature 3

```
1 # Encode class Labels
2 y_train_encoded = label_encoder.fit_transform(y_train3)
3 y_dev_encoded = label_encoder.transform(y_dev3)
4 y_test_encoded = label_encoder.transform(y_test3)
5
6 #Calculate prior, mean and variance of trainset
7 prior, means, var = train_naive_bayes(X_train3, y_train_encoded)
8
9 # Make predictions on the development set
10 dev_predictions = predict_naive_bayes(X_dev3, prior, means, var)
11
12 # Evaluate the model on the development set
13 dev_accuracy = accuracy_score(y_dev_encoded, dev_predictions)
14 print(f'Development Set Accuracy3: {dev_accuracy}')
```

Development Set Accuracy3: 0.571875

Figure 14: Evaluating Naïve Bayes from scratch on three features.

From the figure above it is evident that the model with feature 1 gave a higher accuracy than other features when predicted on the dev set. So, since that feature 1 gave a higher accuracy, it was then used to predict the accuracy on the test set.

```
1 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
2
3 y_test_encoded = label_encoder.transform(y_test1)
4 # Evaluating on test set
5 prior, means, var = train_naive_bayes(X_train1, y_train_encoded)
6 test_predictions = predict_naive_bayes(X_test1, prior, means, var)
7
8
9 print('Accuracy:', accuracy_score(y_test_encoded, test_predictions))
10 print('Confusion matrix:', '\n', confusion_matrix(y_test_encoded, test_predictions))
11 print('Classification Report:', '\n', classification_report(y_test_encoded, test_predictions))
```

Accuracy: 0.71375  
Confusion matrix:  
[[254 124]  
 [105 317]]  
Classification Report:

	precision	recall	f1-score	support
0	0.71	0.67	0.69	378
1	0.72	0.75	0.73	422
accuracy			0.71	800
macro avg	0.71	0.71	0.71	800
weighted avg	0.71	0.71	0.71	800

Figure 15: Evaluation of Naïve Bayes from scratch on test set

From the figure above it is evident that the accuracy the model gave for the test set was **0.71**



### 3.4.1.2 Multinomial Naïve Bayes<sup>11</sup>

The classifier is used for classifying data with discrete features. This is an inbuilt Naïve Bayes Classifier.

Feature 1

```
# Training a machine Learning model (e.g., Naive Bayes) - for Feature 1
1 model_1 = MultinomialNB()
2 model_1.fit(X_train1, y_train1)
3 model_1.predict(X_dev1)
4 # Evaluate the model's performance on the dev data
5 accuracy1 = model_1.score(X_dev1, y_dev1)
6 print(accuracy1)
```

0.765625

Feature 2

```
# Training a machine Learning model (e.g., Naive Bayes) - for Feature 2
1 model_2 = MultinomialNB()
2 model_2.fit(X_train2, y_train2)
3 model_2.predict(X_dev2)
4 # Evaluating the model's performance on the dev data
5 accuracy2 = model_2.score(X_dev2, y_dev2)
6 print(accuracy2)
```

0.784375

Feature 3

```
# Training a machine Learning model (e.g., Naive Bayes) - for Feature 3
1 model_3 = MultinomialNB()
2 model_3.fit(X_train3, y_train3)
3 # Evaluating the model's performance on the dev data
4 model_3.predict(X_dev3)
5 accuracy3 = model_3.score(X_dev3, y_dev3)
6 print(accuracy3)
```

0.7859375

Figure 16: Evaluating Multinomial Naïve Bayes on three features.

Experimenting the three features on the MultinomialNB model gave an accuracy shown in the Figure above and from that it is evident that feature 3 gave a higher accuracy when predicted on the dev set. So, since that model\_3 with feature 3 gave a higher accuracy, it was then used to predict the accuracy on the test set.

```
1 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
2
3 # Evaluating on test set
4 predictions = model_3.predict(X_test3)
5
6 print('Accuracy:', accuracy_score(y_test3, predictions))
7 print('Confusion matrix:', '\n', confusion_matrix(y_test3, predictions))
8 print('Classification Report:', '\n', classification_report(y_test3, predictions))
```

Accuracy: 0.77875  
Confusion matrix:  
[[308 70]  
 [107 315]]  
Classification Report:

	precision	recall	f1-score	support
0	0.74	0.81	0.78	378
1	0.82	0.75	0.78	422
accuracy			0.78	800
macro avg	0.78	0.78	0.78	800
weighted avg	0.78	0.78	0.78	800

Figure 17: Evaluation of Multinomial NB on test set

From the figure above it is evident that the accuracy the model gave for the test set was **0.78**

<sup>11</sup> [https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html)



## 3.4.2 Logistic Regression and Support Vector Machines (SVMs)

### 3.4.2.1 Logistic Regression<sup>12</sup>

Logistic regression is a classification model which is used for linear or binary classification problems. It uses a Logistic function to get the binary outcomes. It does not require a linear relationship between the input and output variables. This is because it applies nonlinear log transformation to the odds ratio, allowing the complex relationship between the variables [34].

$$\text{Logistic Function} = \frac{1}{1 + e^{-x}}$$

These are the default hyperparameters which are set for Logistic Regression: (penalty="l2", \*, dual=False, tol=1e-4, C=1.0, fit\_intercept=True, intercept\_scaling=1, class\_weight=None, random\_state=None, solver="lbfgs", max\_iter=100, multi\_class="auto", verbose=0, warm\_start=False, n\_jobs=None, l1\_ratio=None) [35]

```
Feature 1
M 1 log_reg_1 = LogisticRegression(penalty='l2', dual=False, tol=0.0001,
2                               C=1.0, fit_intercept=True,
3                               intercept_scaling=1,
4                               class_weight=None, random_state=None,
5                               solver='lbfgs', max_iter=100,
6                               multi_class='auto', verbose=0,
7                               warm_start=False, n_jobs=None, l1_ratio=None)
8 log_reg_1.fit(X_train1, y_train1)
9 log_reg_1.predict(X_dev1)
10 score1 = log_reg_1.score(X_dev1, y_dev1)
11 print(f"Accuracy for Dev1: {score1}")
Accuracy for Dev1: 0.6984375

Feature 2
M 1 log_reg_2 = LogisticRegression(penalty='l2', dual=False, tol=0.0001,
2                               C=1.0, fit_intercept=True,
3                               intercept_scaling=1,
4                               class_weight=None, random_state=None,
5                               solver='lbfgs', max_iter=100,
6                               multi_class='auto', verbose=0,
7                               warm_start=False, n_jobs=None, l1_ratio=None)
8 log_reg_2.fit(X_train2, y_train2)
9 log_reg_2.predict(X_dev2)
10 score2 = log_reg_2.score(X_dev2, y_dev2)
11 print(f"Accuracy for Dev2: {score2}")
Accuracy for Dev2: 0.8703125

Feature 3
M 1 log_reg_3 = LogisticRegression(penalty='l2', dual=False, tol=0.0001,
2                               C=1.0, fit_intercept=True,
3                               intercept_scaling=1,
4                               class_weight=None, random_state=None,
5                               solver='lbfgs', max_iter=100,
6                               multi_class='auto', verbose=0,
7                               warm_start=False, n_jobs=None, l1_ratio=None)
8 log_reg_3.fit(X_train3, y_train3)
9 log_reg_3.predict(X_dev3)
10 score3 = log_reg_3.score(X_dev3, y_dev3)
11 print(f"Accuracy for Dev3: {score3}")
Accuracy for Dev3: 0.8546875
```

Figure 18: Evaluating Logistic Regression on three features.

Experimenting the three features on the Logistic Regression model gave an accuracy shown in the Figure above and from that it is evident that feature 2 gave a higher accuracy when predicted on the dev set. So, since the model (log\_reg\_2) with feature 2 gave a higher accuracy, it was then used to predict the accuracy on the test set.

<sup>12</sup> [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)

```

1 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
2
3 # Evaluating on test set
4 predictions = log_reg_2.predict(X_test2)
5
6 print('Accuracy:', accuracy_score(y_test2, predictions))
7 print('Confusion matrix:', '\n', confusion_matrix(y_test2, predictions))
8 print('Classification Report:', '\n', classification_report(y_test2, predictions))

```

Accuracy: 0.84125  
Confusion matrix:  
[[314 64]  
[ 63 359]]

		precision	recall	f1-score	support
	0	0.83	0.83	0.83	378
	1	0.85	0.85	0.85	422
	accuracy			0.84	800
	macro avg	0.84	0.84	0.84	800
	weighted avg	0.84	0.84	0.84	800

Figure 19: Evaluation of Logistic Regression on test set

From the figure above it is evident that the accuracy the model gave for the test set was **0.84**

### 3.4.2.2 SVM (Support Vector Machines)<sup>13</sup>

A supervised learning method which is used for classification, regression, and outliers' detections tasks. Effective in high dimensional tasks where number of dimensions is greater than the number of samples. SVC (Support Vector Classification) is the one which is responsible for binary and multi-class classification tasks.

These are the default hyperparameters which are set for SVC:  
(\*, C=1.0, kernel='rbf', degree=3, gamma="scale", coef0=0.0, shrinking=True, probability=False, tol=1e-3, cache\_size=200, class\_weight=None, verbose=False, max\_iter=-1, decision\_function\_shape="ovr", break\_ties=False, random\_state=None.)  
[36,37].

Feature 1	Feature 2	Feature 3
<pre> 1 svm_clf_1 = SVC(C=1.0, kernel='rbf', 2 degree=3, gamma='scale', 3 coef0=0.0, shrinking=True, 4 probability=False, tol=0.001, 5 cache_size=200, class_weight=None, 6 verbose=False, max_iter=-1, 7 decision_function_shape='ovr', 8 break_ties=False, random_state=None) 9 10 svm_clf_1.fit(X_train1, y_train1) 11 svm_clf_1.predict(X_dev1) 12 score1 = svm_clf_1.score(X_dev1, y_dev1) 13 print(f"Accuracy for Dev1: {score1}") </pre> <p>Accuracy for Dev1: 0.7453125</p>	<pre> 1 svm_clf_2 = SVC(C=1.0, kernel='rbf', 2 degree=3, gamma='scale', 3 coef0=0.0, shrinking=True, 4 probability=False, tol=0.001, 5 cache_size=200, class_weight=None, 6 verbose=False, max_iter=-1, 7 decision_function_shape='ovr', 8 break_ties=False, random_state=None) 9 10 svm_clf_2.fit(X_train2, y_train2) 11 svm_clf_2.predict(X_dev2) 12 score2 = svm_clf_2.score(X_dev2, y_dev2) 13 print(f"Accuracy for Dev2: {score2}") </pre> <p>Accuracy for Dev2: 0.8375</p>	<pre> 1 svm_clf_3 = SVC(C=1.0, kernel='rbf', 2 degree=3, gamma='scale', 3 coef0=0.0, shrinking=True, 4 probability=False, tol=0.001, 5 cache_size=200, class_weight=None, 6 verbose=False, max_iter=-1, 7 decision_function_shape='ovr', 8 break_ties=False, random_state=None) 9 10 svm_clf_3.fit(X_train3, y_train3) 11 svm_clf_3.predict(X_dev3) 12 score3 = svm_clf_3.score(X_dev3, y_dev3) 13 print(f"Accuracy for Dev3: {score3}") </pre> <p>Accuracy for Dev3: 0.8375</p>

Figure 20: Evaluating SVC on three features.

Experimenting the three features on the SVC model gave an accuracy shown in the Figure above and from that it is evident that feature 2 and feature 3 gave same accuracy when predicted on the dev set. So, since both models gave the same accuracy, model (svm\_clf\_3) with feature 3 was taken to check the accuracy on the test set.

<sup>13</sup> <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>

```

1 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
2
3 # Now apply those above metrics to evaluate your model
4 predictions = svm_clf_3.predict(X_test3)
5
6 print('Accuracy:', accuracy_score(y_test3, predictions))
7 print('Confusion matrix:', '\n', confusion_matrix(y_test3, predictions))
8 print('Classification Report:', '\n', classification_report(y_test3, predictions))

```

Accuracy: 0.805  
Confusion matrix:  
[[265 113]  
[ 43 379]]

Classification Report:

	precision	recall	f1-score	support
0	0.86	0.70	0.77	378
1	0.77	0.90	0.83	422
accuracy			0.81	800
macro avg	0.82	0.80	0.80	800
weighted avg	0.81	0.81	0.80	800

Figure 21: Evaluation of SVC on test set

From the figure above it is evident that the accuracy the model gave for the test set was **0.81**

### 3.4.2.2 Hyperparameter Optimisation<sup>14</sup>

Some of the Hyperparameters of both Logistic Regression and SVM was changed and tested to see the performance on the model which gave a higher accuracy when default parameters were used. From that the model with higher accuracy after the hyperparameters optimisation was tested on the test dataset to see the performance.

**Table 1:** Hyperparameter Optimisation of Logistic Regression Classifier

Parameter	log_reg_4	log_reg_5	log_reg_6	log_reg_7	log_reg_8
<b>Penalty</b>	L1	L2	L2	L2	L2
<b>Dual</b>	False	False	False	False	False
<b>Tol</b>	0.0001	0.0001	0.0001	0.0001	0.0001
<b>C</b>	1.0	1.0	1.0	1.0	1.0
<b>Fit Intercept</b>	True	True	True	True	True
<b>Intercept Scaling</b>	1	1	1	1	1
<b>Class Weight</b>	Balanced	Balanced	Dict	Dict	Balanced
<b>Random State</b>	None	None	None	None	42
<b>Solver</b>	liblinear	lbfgs	lbfgs	liblinear	sag
<b>Max Iter</b>	1000	1000	100	10000	100
<b>Multi-Class</b>	Auto	Multinomial	Auto	Auto	Auto
<b>Verbose</b>	5	0	0	10	10
<b>Warm Start</b>	False	False	False	False	False
<b>N Jobs</b>	None	None	None	None	None
<b>L1 Ratio</b>	None	None	None	None	None
<b>Accuracy</b>	<b>0.846875</b>	<b>0.8671875</b>	<b>0.8703125</b>	<b>0.875</b>	<b>0.88125</b>

<sup>14</sup> [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)  
<https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>

```

1 #penalty=l2, class_weight='balanced',max_iter=100, solver='sag',random_state=42
2 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
3
4 # Now apply those above metrics to evaluate your model
5 predictions = log_reg_8.predict(X_test2)
6
7 print('Accuracy:',accuracy_score(y_test2,predictions))
8 print('Confusion matrix:','\n',confusion_matrix(y_test2,predictions))
9 print('Classification Report:','\n',classification_report(y_test2,predictions))

```

Accuracy: 0.83375  
 Confusion matrix:  
 [[310 68]  
 [ 65 357]]  
 Classification Report:

	precision	recall	f1-score	support
0	0.83	0.82	0.82	378
1	0.84	0.85	0.84	422
accuracy			0.83	800
macro avg	0.83	0.83	0.83	800
weighted avg	0.83	0.83	0.83	800

Figure 22: Evaluation on test set of Logistic Regression after hyperparameter Optimisation

- The evaluation of the test set using the optimal hyperparameters for log\_reg\_8 resulted in an accuracy of **0.83**.
- While running Logistic Regression with default parameters on feature 1, a warning of 'Max Iterations reached' was encountered. To address this, the Max Iterations parameter was adjusted from 100 to 1000 and 10000.

```

STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear\_model.html#logistic-regression
n_iter_i = _check_optimize_result(

```

Figure 23: Maximum Iteration reached Warning.

- Table 1 highlights the performance of the data with Logistic Regression, particularly when various hyperparameters are modified. Testing using different solvers, 'liblinear,' 'lbfgs,' and 'sag'—revealed that the 'sag' solver with a random state of 42 and with other default parameters produced the highest accuracy.
- Additional parameters, such as Class Weight (dict or Balanced), Max Iterations (1000, 10000), and multi-class (Auto, Multinomial), were explored by having different solvers to check which gave a better accuracy.
- These are the few hyperparameters which have been randomly changed and tested, a lot of combinations can be further tested.

**Table 2:** Hyperparameter Optimisation of SVM (SVC) Classifier

Parameter	svm_clf_4	svm_clf_5	svm_clf_6	svm_clf_7	svm_clf_8
<b>Kernel</b>	Poly	Sigmoid	RBF	Sigmoid	Poly
<b>Degree</b>	3	3	3	3	3
<b>Gamma</b>	Auto	Scale	Auto	Scale	Scale
<b>Coef0</b>	0.1	0.1	0.1	0.1	0.1
<b>Shrinking</b>	True	True	True	True	True
<b>Probability</b>	False	False	False	False	False
<b>Tol</b>	0.001	0.001	0.001	0.001	0.001
<b>Cache Size</b>	200	200	200	200	200

<b>Class Weight</b>	Balanced	None	Balanced	Balanced	Balanced
<b>Verbose</b>	False	False	False	False	False
<b>Max Iter</b>	-1	-1	-1	-1	-1
<b>Decision Function Shape</b>	OVO	OVO	OVO	OVR	OVR
<b>Break Ties</b>	False	False	False	False	False
<b>Random State</b>	42	42	42	42	None
<b>Accuracy</b>	<b>0.4859375</b>	<b>0.85625</b>	<b>0.515625</b>	<b>0.853125</b>	<b>0.540625</b>

```

1 #kernel='sigmoid', gamma='scale', coef0 = 0.1,
2 #class_weight= None, decision_function_shape='ovo',
3 #random_state=42
4 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
5
6 # Now apply those above metrics to evaluate your model
7 predictions = svm_clf_5.predict(X_test3)
8
9 print('Accuracy:', accuracy_score(y_test3, predictions))
10 print('Confusion matrix:', '\n', confusion_matrix(y_test3, predictions))
11 print('Classification Report:', '\n', classification_report(y_test3, predictions))

```

Accuracy: 0.83625  
 Confusion matrix:  
 [[305 73]  
 [ 58 364]]  
 Classification Report:  

	precision	recall	f1-score	support
0	0.84	0.81	0.82	378
1	0.83	0.86	0.85	422
accuracy			0.84	800
macro avg	0.84	0.83	0.84	800
weighted avg	0.84	0.84	0.84	800

Figure 24: Evaluation on test set of SVC after hyperparameter Optimisation

- The test set evaluation using the best hyperparameters for svm\_clf\_5 gave an accuracy of **0.84**.
- Analysis from Table 2 indicates that the SVM model performs well on this data when the kernel is set to 'Sigmoid' with a random state of 42. Conversely, it exhibits suboptimal performance with the 'Poly' kernel.
- The default kernel for SVC is 'RBF,' exploration of other hyperparameters such as Gamma (from 'Scale' to 'Auto'), Class Weight (from 'None' to 'Balanced'), and Decision Function shape (from 'OVR' to 'OVO') also gave suboptimal performance.
- The main hyperparameters that were tested are different kernels (Poly, Sigmoid, and RBF). Depending on the chosen kernel, additional hyperparameters were adjusted. Notably, Gamma was assessed using both 'Auto' and 'Scale.'
- These are the few hyperparameters which have been randomly changed and tested, a lot of combinations can be further tested.

## 3.5 Testing on Deep Learning Algorithm

### 3.5.1 BERT

**BERT<sup>15</sup> (Bidirectional Encoder Representation for Transformers)** is a Machine Learning transformer model which is used for many different applications and one which is sentiment

<sup>15</sup> [https://huggingface.co/transformers/v3.2.0/custom\\_datasets.html](https://huggingface.co/transformers/v3.2.0/custom_datasets.html)  
[https://huggingface.co/docs/transformers/tasks/sequence\\_classification](https://huggingface.co/docs/transformers/tasks/sequence_classification)  
[https://github.com/dipanjanS/deep\\_transfer\\_learning\\_nlp\\_dhs2019/tree/master/notebooks](https://github.com/dipanjanS/deep_transfer_learning_nlp_dhs2019/tree/master/notebooks)

analysis. This method is useful when there are large amounts of data, and this transform model brings parallelization which makes it easier to train large amounts of data [38].  
List of Libraries imported to run the BERT model:

```
1 #Libraries
2
3 import numpy as np
4 import pandas as pd
5 import nltk
6 import string
7 import os
8 import evaluate
9 import torch
10
11
12 from sklearn.model_selection import train_test_split
13 from transformers import DistilBertTokenizerFast
```

Figure 25: List of Libraries

Data Split<sup>16</sup> : Splitting dataset into Train, Development, and Test Set and Formatting them into CSV files. The dataset for BERT was divided into train, development, and test splits and all of them having same number of splits as tested with other models. That is the whole data is first divided into train and test set with train set having 80% of data and test set having 20%. Then the train set is then further divided into train and development set in which train set having 80% of data and development set having 20% of the data.

The data split into train, test and dev sets are then formatted into .csv files for further processing.

```
1 # Taking 'reviews' and 'Labels' columns from the whole dataset
2 data = {'Reviews': df['reviews'], 'Labels': df['labels']}
3 df_bert = pd.DataFrame(data)
4
5
6
7 # Split data into train, dev, and test sets
8 X_train, X_test, y_train, y_test = train_test_split(df_bert['Reviews'], df_bert['Labels'],
9                                                    test_size=0.2, random_state=42)
10 X_train, X_dev, y_train, y_dev = train_test_split(X_train, y_train, test_size=0.2,
11                                                    train_size=0.8, random_state=42)
12
13 train_df = pd.DataFrame({'Reviews': X_train, 'Labels': y_train})
14 dev_df = pd.DataFrame({'Reviews': X_dev, 'Labels': y_dev})
15 test_df = pd.DataFrame({'Reviews': X_test, 'Labels': y_test})
16
17 # Save to CSV files
18 train_df.to_csv('train.csv', index=False)
19 dev_df.to_csv('dev.csv', index=False)
20 test_df.to_csv('test.csv', index=False)
```

Figure 26: Splitting data into train, test and dev sets and converting them into csv files.

Two pre-trained versions of BERT were used for this sentiment analysis. BERT- base cased and BERT- base uncased. In uncased version all the words are converted to lowercase and in cased version all the uppercase and lowercase words are kept as it is [39].

Tokenization was done using DistilBertTokenizerFast for cased and uncased.

<sup>16</sup> <https://stackoverflow.com/questions/38250710/how-to-split-data-into-3-sets-train-validation-and-test>

```

1 from transformers import DistilBertTokenizerFast
2 tokenizer = DistilBertTokenizerFast.from_pretrained('distilbert-base-cased')

1 train_encodings = tokenizer(train_df['Reviews'].tolist(), truncation=True, padding=True)
2 dev_encodings = tokenizer(dev_df['Reviews'].tolist(), truncation=True, padding=True)
3 test_encodings = tokenizer(test_df['Reviews'].tolist(), truncation=True, padding=True)

```

Figure 27: Tokenization of based cased.

```

1 #Tokenizer for Uncased version of the Base model
2 from transformers import DistilBertTokenizerFast
3 tokenizer1 = DistilBertTokenizerFast.from_pretrained('distilbert-base-uncased')

1 train_encodings1 = tokenizer1(train_df['Reviews'].tolist(), truncation=True, padding=True)
2 dev_encodings1 = tokenizer1(dev_df['Reviews'].tolist(), truncation=True, padding=True)
3 test_encodings1 = tokenizer1(test_df['Reviews'].tolist(), truncation=True, padding=True)

```

Figure 28: Tokenization of based uncased.

A class called IMDbDataset was created which takes all the labels and the encoded reviews and processed so that PyTorch datasets can be easily used with a PyTorch DataLoader during the training process [46].

```

1 import torch
2
3 class IMDbDataset(torch.utils.data.Dataset):
4     def __init__(self, encodings, labels):
5         self.encodings = encodings
6         self.labels = labels
7
8     def __getitem__(self, idx): #Asked chatGPT
9         if idx < len(self.labels):
10             item = {key: torch.tensor(val[idx]) for key, val in self.encodings.items()}
11             item['labels'] = torch.tensor(self.labels.iloc[idx]) # Use iloc to access the Label by index
12             return item
13         else:
14             # Handle the case where the index is out of range
15             return None # or raise an exception, depending on your needs
16
17     def __len__(self):
18         return len(self.labels)
19
20
21 train_dataset = IMDbDataset(train_encodings, train_df['Labels'])
22 dev_dataset = IMDbDataset(dev_encodings, dev_df['Labels'])
23 test_dataset = IMDbDataset(test_encodings, test_df['Labels'])
24
25

```

Figure 29: Labels and reviews proceed so it can be used with PyTorch.

Then the transformed dataset is then given into the transform model for Sequence Classification (DistilBertSequenceClassification).

```

from transformers import DistilBertForSequenceClassification, Trainer, TrainingArguments

```

Figure 30: Loading the transformer model for training.

These are the arguments(hyperparameters) which are used to train the model. For the model to run faster, the number of epochs was reduced to 1 and learning rate was also reduced to a smaller value 1e-5. Then the model is trained on the trainset and then evaluated on the development set. That is Fine-Tuning the BERT model on the dataset provided.

```

from transformers import DistilBertForSequenceClassification, Trainer, TrainingArguments

training_args = TrainingArguments(
    output_dir='./results',
    num_train_epochs=1,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=64,
    learning_rate=1e-5,
    warmup_steps=10,
    weight_decay=0.01,
    logging_dir='./logs',
    logging_steps=10,
)

model = DistilBertForSequenceClassification.from_pretrained("distilbert-base-cased")

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=dev_dataset,
    compute_metrics=compute_metrics
)

trainer.train()

```

Figure 31: Setting parameters and fine-tuning the model.

This model is used for cased and uncased versions.

```

model = DistilBertForSequenceClassification.from_pretrained("distilbert-base-uncased")

model = DistilBertForSequenceClassification.from_pretrained("distilbert-base-cased")

```

Figure 32: Model for cased and uncased

```

1 #To evaluate the performace of the transformer
2 def compute_metrics(eval_pred):
3     predictions, labels = eval_pred
4     predictions = np.argmax(predictions, axis=1)
5     return accuracy.compute(predictions=predictions, references=labels)

```

```

1 accuracy = evaluate.load("accuracy")

```

Figure 33: Evaluating<sup>17</sup> the performance of the model.

The Accuracy of the model was evaluated first on the Dev dataset and from the Figures below BERT-base cased version gave an accuracy of 0.865 and BERT-base uncased version gave an accuracy of 0.875.

```

1 trainer.evaluate()

```

```

{'eval_loss': 0.36846449971199036,
 'eval_accuracy': 0.865625,
 'eval_runtime': 482.9082,
 'eval_samples_per_second': 1.325,
 'eval_steps_per_second': 0.021,
 'epoch': 1.0}

```

Figure 34: Accuracy on the dev set of BERT cased version.

<sup>17</sup> <https://www.kaggle.com/code/pritishmishra/text-classification-with-distilbert-92-accuracy>  
[https://huggingface.co/docs/transformers/tasks/sequence\\_classification](https://huggingface.co/docs/transformers/tasks/sequence_classification)



```

1 trainer.evaluate()

{'eval_loss': 0.3671916127204895,
 'eval_accuracy': 0.875,
 'eval_runtime': 480.8928,
 'eval_samples_per_second': 1.331,
 'eval_steps_per_second': 0.021,
 'epoch': 1.0}

```

*Figure 35: Accuracy on the dev set of BERT uncased version.*

Then Accuracy on the test dataset was checked to see how the model was performing. The cased version gave an accuracy of 0.845 for the test set and uncased version gave an accuracy of 0.878.

```

1 pred_test = trainer.predict(test_dataset)

1 pred_test.metrics

{'test_loss': 0.38908982276916504,
 'test_accuracy': 0.845,
 'test_runtime': 621.418,
 'test_samples_per_second': 1.287,
 'test_steps_per_second': 0.021}

```

*Figure 36: Accuracy on the test set of BERT cased version*

```

1 pred_test1 = trainer.predict(test_dataset1)

1 pred_test1.metrics

{'test_loss': 0.367340087890625,
 'test_accuracy': 0.87875,
 'test_runtime': 608.4871,
 'test_samples_per_second': 1.315,
 'test_steps_per_second': 0.021}

```

*Figure 37: Accuracy on the test set of BERT uncased version*

## 4. Discussion

**Table 3:** Results of each model on the Test set

Method	Accuracy
Naïve Bayes (from scratch)	0.71
Naïve Bayes (scikit-learn)	0.78
Logistic Regression	0.84
SVM	0.81
Logistic Regression (with Hyperparameter Optimisation)	0.83
SVM (with Hyperparameter Optimisation)	0.84
BERT (Cased)	0.845
BERT (Uncased)	0.87875

Table 3 shows the overall result of the test data on each of the machine learning models. BERT-base uncased version gave the highest accuracy compared to the other models with accuracy **~0.88**.

Naïve Bayes from scratch gave a lower accuracy of 0.71 this maybe because of the way the implementation was done and by changing the way the distribution was calculated, that is instead of using Gaussian Distribution maybe calculating the distribution using Multinomial or Bernoulli [42] would have given a higher accuracy. Multinomial Naïve bayes gave an accuracy of 0.78 for test set. Which again is not that high compared to accuracies given by other models. So overall I feel Naïve Bayes is not a good model for classifying this dataset.

Next algorithm the data was tested on was Logistic Regression and it gave an accuracy of 0.84 with default hyperparameters and when the hyperparameters of the model was changed to calculate the performance of the model it gave an accuracy of 0.83. From testing it was clear that this data works well with Logistic Regression model. Even after hyperparameter Optimisation the model was giving a higher accuracy but the accuracy on the test was slightly lower.

Then it was tested on SVM (SVC) where again it performed almost like Logistic Regression and gave an accuracy of 0.81 with default parameter and 0.84 after Hyperparameter Optimisation. So, again this model is also a good choice for this dataset.

Some reasons why Logistic Regression and SVC outperformed Naïve Bayes is because these algorithms try to capture more complex information from the data compared to Naïve Bayes where it checks for the conditional probability of the feature of the give class to make assumptions and then classify the model. Both Logistic Regression and SVC can handle non-linear relationship between the transformed data and Naïve Bayes being linear may not be able to handle it. For Naïve Bayes there is not Hyperparameter Optimisation so no changes can be made to the model to improve the performance [43].

The IMDB reviews dataset was then applied onto the BERT Transformer model. This model outperformed all the other classification models. BERT-base uncased version gave the highest accuracy when compared to all other models. This is because BERT having this transformer architecture allows for parallelization during training so as result it will be able to capture more long-range dependencies of the texts [45]. Since the data is being fine-

tuned on the pre-trained model (a model which has been trained on a large corpus of data), the model will be able to capture any intricate information and patterns from the data and can more accurately make the classification [47].

## 5. Conclusion and Future Work

In this project the IMDB Reviews dataset was tested on a few Machine Learning algorithms and with a Deep Learning algorithm (BERT).

Overall, the best performing model for this IMDB Reviews dataset is BERT-base uncased version. This is because BERT, being a pretrained model that is it has been trained on a large corpus of data the model is able to capture any form of patterns from the data after fine-tuning. Fine-tuning on the pre-trained model, it can adapt to the data provided and make changes accordingly for sentiment analysis and BERT model is able to handle any complex tasks. The accuracy of the BERT model can further be improved by proper Hyperparameter tuning and using different Tokenization algorithms [44].

To improve the results in the future I would consider using built-in functions for TF-IDF which will give better accuracy for the features generated. Trying out many more feature generation and feature selection techniques. Testing with more hyperparameters for Logistic Regression and SVM. For the BERT model increasing the learning rate and epochs. For Naïve Bayes from scratch try to implement most of the code by myself in the future to see the changes in the accuracy. Due to time constraints these could not be done and tested.

## References

- I acknowledge that this work is my own, and I have used ChatGPT 3.5 (OpenAI, <https://chat.openai.com/>) in few places for code development and to proofread the final draft of my report.
- 1. Learning Word Vectors for Sentiment Analysis, Maas, Andrew L. and Daly, Raymond E. and Pham, Peter T. and Huang, Dan and Ng, Andrew Y. and Potts, Christopher, June 2011, online - <http://www.aclweb.org/anthology/P11-1015> , <https://ai.stanford.edu/~amaas/data/sentiment/>
- 2. MonkeyLearn, Sentiment Analysis: A Definitive Guide, online: <https://monkeylearn.com/sentiment-analysis/>
- 3. THE APP SOLUTIONS, Why Business Applies Sentiment Analysis? 5 Successful Examples, 2023, online: <https://theappsolutions.com/blog/development/sentiment-analysis-for-business/>
- 4. AIMultiple, Top 5 Sentiment Analysis Challenges and Solutions in 2023, online: <https://research.aimultiple.com/sentiment-analysis-challenges/>
- 5. Dovetail Editorial Team, 2023, Understanding sentiment analysis using machine learning online: <https://dovetail.com/customer-research/sentiment-analysis-using-machine-learning/#:~:text=Four%20machine%20learning%20techniques%20for,reinforcement%20and%20semi%2Dsupervised%20learning.>

6. DataRobot, 2019, Using Machine Learning for Sentiment Analysis: a Deep Dive, online: <https://www.datarobot.com/blog/using-machine-learning-for-sentiment-analysis-a-deep-dive/#:~:text=RNNs%20are%20probably%20the%20most,piece%20of%20text%20is%20ingested>.
7. Yulia Gavrilova, 2023, Transformers in ML: What They Are and How They Work, online: <https://serokell.io/blog/transformers-in-ml>
8. C.J. Hutto, Eric Gilbert, January 2015, VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text.  
Online:[https://www.researchgate.net/publication/275828927\\_VADER\\_A\\_Parsimoniou\\_s\\_Rule-based\\_Model\\_for\\_Sentiment\\_Analysis\\_of\\_Social\\_Media\\_Text](https://www.researchgate.net/publication/275828927_VADER_A_Parsimoniou_s_Rule-based_Model_for_Sentiment_Analysis_of_Social_Media_Text)
9. Kavya Suppala, Narasinga Rao, June 2019, Sentiment Analysis using Naïve Bayes Classifier, online: <https://www.ijitee.org/wp-content/uploads/papers/v8i8/H6330068819.pdf>
10. Christopher D. Manning Prabhakar Raghavan Hinrich Schütze, 2009, An Introduction to Information Retrieval, online: <https://nlp.stanford.edu/IR-book/pdf/irbookonlinereading.pdf>
11. Tony Mullen and Nigel Collier, Sentiment analysis using support vector machines with diverse information sources, online: <https://aclanthology.org/W04-3253.pdf>
12. Corinna Cortes, Vladimir Vapnik, 1995, Support- Vector Networks, online: <https://link.springer.com/article/10.1023/A:1022627411411>
13. Leo Breiman, 2001, Random Forests, online: <https://www.stat.berkeley.edu/~breiman/randomforest2001.pdf>
14. Shoffan Saifullah<sup>1,a</sup>, Yuli Fauziah<sup>1,b</sup>, Agus Sasmito Aribowo<sup>1,2,c</sup>, Comparison of Machine Learning for Sentiment Analysis in Detecting Anxiety Based on Social Media Data, online: [https://arxiv.org/ftp/arxiv/papers/2101/2101.06353.pdf#:~:text=Random%20Forest%20\(RF\)%20is%20a,the%20sample%20data%20they%20have](https://arxiv.org/ftp/arxiv/papers/2101/2101.06353.pdf#:~:text=Random%20Forest%20(RF)%20is%20a,the%20sample%20data%20they%20have).
15. Sepp Hochreiter, Jürgen Schmidhuber, 1997, Long Short-Term Memory, online: <https://www.bioinf.jku.at/publications/older/2604.pdf>
16. Robin M. Schmidt, November 2019, Recurrent Neural Networks (RNNs): A gentle Introduction and Overview, online: <https://arxiv.org/pdf/1912.05911.pdf>
17. Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, January 2019, BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, online: <https://arxiv.org/pdf/1810.04805v2.pdf>
18. Bruno Stecanella, June 2017, MonkeyLearn, Support Vector Machines (SVM) Algorithm Explained, online: <https://monkeylearn.com/blog/introduction-to-support-vector-machines-svm/>
19. Intellipaat Software Solutions, What is LSTM? Introduction to Long Short-Term Memory, online: <https://intellipaat.com/blog/what-is-lstm/>
20. IBM, What are Recurrent Neural Networks?, online: [https://www.ibm.com/topics/recurrent-neural-networks#:~:text=A%20recurrent%20neural%20network%20\(RNN,data%20or%20time%20series%20data](https://www.ibm.com/topics/recurrent-neural-networks#:~:text=A%20recurrent%20neural%20network%20(RNN,data%20or%20time%20series%20data).
21. Ben Lutkevich, BERT Language mode, online: <https://www.techtarget.com/searchenterpriseai/definition/BERT-language->



37. Sklearn.svm.SVC, online: <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC>
38. Britney Muller, BERT 101 State of the Art NLP Model Explained, March 2022, online: <https://huggingface.co/blog/bert-101>
39. K. Kyritsis, N. Spatiotis, I. Perokos, M. Paraskevas, A Comparative Performance Evaluation of Algorithms for the Analysis and Recognition of Emotional Content, DOI: 10.5772/intechopen.112627, July 2023, online: <https://www.intechopen.com/online-first/87923>
40. Yellowbrick, Classification Report, online: [https://www.scikit-yb.org/en/latest/api/classifier/classification\\_report.html#:~:text=The%20classification%20report%20visualizer%20displays,with%20a%20color%2Dcoded%20heatmap.](https://www.scikit-yb.org/en/latest/api/classifier/classification_report.html#:~:text=The%20classification%20report%20visualizer%20displays,with%20a%20color%2Dcoded%20heatmap.)
41. Sarang Narkhede, Understanding Confusion Matrix, 2018, online: <https://towardsdatascience.com/understanding-confusion-matrix-a9ad42dcfd62>
42. Nagesh Singh Chauhan, Naïve Bayes Algorithm: Everything You Need To Know, April 2022, Online: <https://www.kdnuggets.com/2020/06/naive-bayes-algorithm-everything.html>
43. P. Golpour , M. Ghayour-Mobarhan, A. Sak., H. Esmaily, A. Taghipour, M. Tajfard, H. Ghazizadeh, M. Moohebat, G.A. Ferns, 2020, Comparison of Support Vector Machine, Naïve Bayes and Logistic Regression for Assessing the Necessity for Coronary Angiography. International journal of environmental research and public health, 17(18), 6449. Online: <https://doi.org/10.3390/ijerph17186449>
44. S. Alaparhi, M. Mishra, BERT: a sentiment analysis odyssey, J Market Anal 9, 118–126, 2021, online: <https://doi.org/10.1057/s41270-021-00109-8>
45. Edoardo Guerriero, Why does the transformer do better than RNN and LSTM in long-range context dependencies?, Stack Exchange, online: <https://ai.stackexchange.com/questions/20075/why-does-the-transformer-do-better-than-rnn-and-lstm-in-long-range-context-depen>
46. Hugging Face, Fine-Tuning with custom Datasets, online: [https://huggingface.co/transformers/v3.2.0/custom\\_datasets.html](https://huggingface.co/transformers/v3.2.0/custom_datasets.html)
47. Nikolaj Buhl, Training vs. Fine-tuning: What is the Difference, Nov 2023, online: <https://encord.com/blog/training-vs-fine-tuning/>