



كلية علوم الحاسب والمعلومات  
قسم هندسة البرمجيات

King Saud University

College of Computer and Information Sciences Department

of Software Engineering

SWE 486 – Cloud Computing and Big Data



# Movie Reviews

## Sentiment Analysis

PHASE 2

---

#	NAME	ID
1	Lina Alkhodhiri	441201109
2	Noura Alsaif	439200832
3	Alanoud Alomar	441201113

GROUP #	4
SUPERVISOR	L. Hailah Almazrui

*Submission Date: August 14, 2022*

## Table of contents

1. Descriptive Analytics .....	1
2. Predictive Analytics .....	4
2.1 Naïve Bayes Model .....	6
3. Discussion questions .....	10
4. Reference.....	11

## Table of figures

Figure 1.....	1
Figure 2.....	1
Figure 3.....	1
Figure 4.....	2
Figure 5.....	2
Figure 6.....	3
Figure 7.....	3
Figure 8.....	4
Figure 9.....	5
Figure 10.....	6
Figure 11.....	7
Figure 12.....	7
Figure 13.....	7
Figure 14.....	8
Figure 15.....	8
Figure 16.....	9

## Phase 2 Model planning and building

### 1. Descriptive Analytics

Descriptive Analysis is the type of data that helps describe, show, or summarize data points in a constructive way such that patterns might emerge that fulfill every condition of the data. It is one of the most important steps for conducting statistical data analysis.

- **import libraries**

```
import pandas as pd
import numpy as np
```

*Figure 1. import libraries*

- **Shape and columns of the data frame**

```
✓ [131] data_df.shape
0s
      (1524, 5)

✓ [133] data_df.columns
0s
      Index(['Text', 'classification', 'clean text', 'tokens', 'text length'], dtype='object')
```

*Figure 2. data frame columns and shape*

The shape of the data frame has 1524 rows and 5 columns.

- **data frame summary**

```
✓ [132] data_df.info(verbose = True)
0s
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1524 entries, 0 to 1523
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Text             1524 non-null   object
1   classification    1519 non-null   object
2   clean text       1524 non-null   object
3   tokens           1524 non-null   object
4   text length      1524 non-null   int64
dtypes: int64(1), object(4)
memory usage: 59.7+ KB
```

*Figure 3. data frame summary*

.info() function is used by pandas to show a brief summary about a DataFrame including the index dtype and columns, non-null values, and memory usage.

- **Statistics**

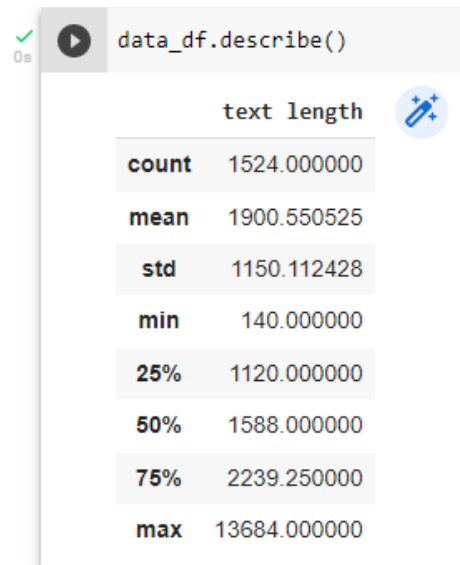


Figure 4. view statistical details

In figure 4 the result after we used describe() method which is helping to view some basic statistical details like percentile, mean, std, min, max of the data frame.

- **Word occurrences**

We had to read a data frame to count the idf for getting term frequency.

```
[160] #libraries for word occurrence()
      from sklearn.feature_extraction.text import TfidfTransformer
      from sklearn.feature_extraction.text import CountVectorizer

[167] #instantiate CountVectorizer()
      countVec = CountVectorizer()

      #generate word counts for the words
      word_count_vector = countVec.fit_transform(data_df['Text'].astype('U'))
      word_count_vector.shape

      (1524, 72745)

[168] #Transform a count matrix to a normalized tf-idf representation
      tfidf_transformer = TfidfTransformer(smooth_idf=True, use_idf=True)
      #idf values
      tfidf_transformer.fit(word_count_vector)

      TfidfTransformer()

[169] #print idf values
      df_idf = pd.DataFrame(tfidf_transformer.idf_, index = countVec.get_feature_names(), columns = ["idf_weights"])
```

Figure 5. import sklearn libraries to count the idf

- **Most frequent terms**



Figure 6. most frequent terms using sort by idf

- **Least frequent terms**

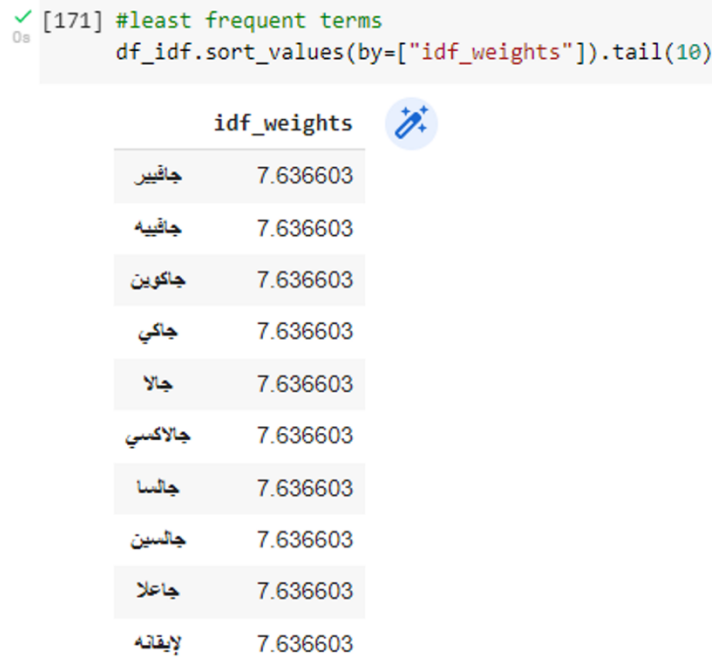


Figure 7. least frequent terms using sort by idf

## 2. Predictive Analytics

In our project, we chose the Naïve Bayes model and used several techniques to evaluate the results. We will discuss the steps of our implementation of the Naïve Bayes but before that, we will mention the advantages of the model and the definition of the technologies that we used:

- **Naïve Bayes Model and its features**

Naïve Bayes algorithm is a supervised learning algorithm used for solving classification problems and it is mainly used in text classification that includes a high-dimensional training dataset. Naïve Bayes Classifier is one of the simplest and most effective Classification algorithms which helps in building fast machine learning models that can make quick predictions. Some popular examples of Naïve Bayes Algorithm are spam filtration, classifying articles, or in our case Sentimental analysis [2].

- **ROC Curve**

An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds. It basically shows two parameters: true positive rate and false positive rate. And by the looks of it, we can determine the best performance by observing how the curve is plotted as demonstrated in the figure below [3].

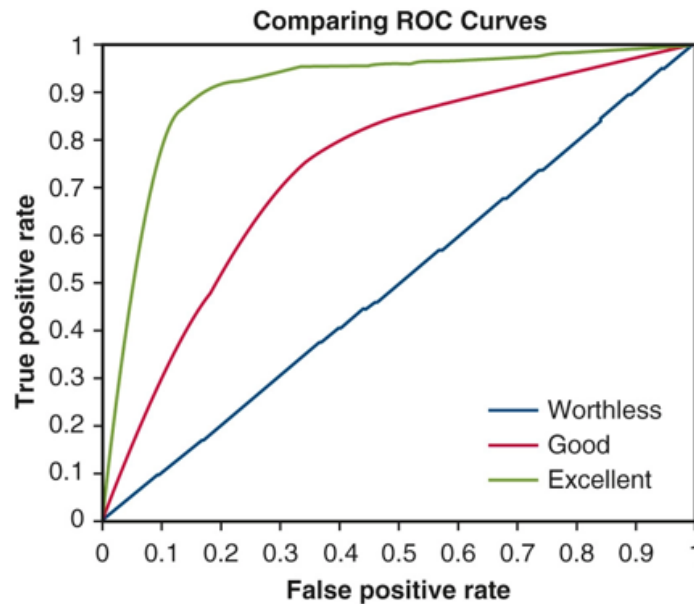


Figure 8. ROC Curve

- **Ten-Folds Cross Validation**

The basic idea of ten folds cross validation, as the name implies, is to divide the data set into 10 folds. First, we will iterate 10 times over the data, each time holding one-fold for testing and the 9 remaining folds will be used for training. And so on until the 10 folds had all been used for testing once and for training 9 times. The accuracy will be the average accuracy of all 10 iterations. Which makes ten folds cross validation results in skill estimates that generally have a lower bias than other methods.

- **Confusion Matrix**

The confusion matrix provides many helpful evaluation measurements. Precision, Recall, ...etc. are all derived from the confusion matrix. Which essentially provides counts of true positive, false positive, true negative and false negative predictions as illustrated in figure 11. Precision - Precision is the ratio of correctly predicted positive observations to the total predicted positive observations. Recall (Sensitivity) - Recall is the ratio of correctly predicted positive observations to all observations in actual class.

## Confusion Matrix

	Actually Positive (1)	Actually Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)

*Figure 9. Confusion Matrix*

## 2.1 Naïve Bayes Model

- **Libraries used**

**NumPy:** It is a Python library that provides a multidimensional array object, various derived objects, and an assortment of routines for fast operations on arrays, including mathematical, logical, shape manipulation, sorting, selecting, I/O, discrete Fourier transforms, basic linear algebra, basic statistical operations, random simulation and much more [4].

**Sklearn:** Is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering, and dimensionality reduction via a consistency interface in Python [5].

**Pandas:** Pandas is an open-source Python package that is most widely used for data science/data analysis and machine learning tasks. It is built on top of another package named Numpy, which provides support for multi-dimensional arrays [6].

**Matplotlib:** Matplotlib is a cross-platform, data visualization and graphical plotting library for Python and its numerical extension NumPy [7].

- **steps of our implementation**

- 1- We have removed the neutral rows since the neutral class is incomprehensible and does not relate to a specific and clear class.

```
# remove the "Neutral" class  
data_df=data_df[data_df['classification'] != "Neutral"]
```

*Figure 10. Remove neutral class*



- To process each text's class label, we converted the values from "Positive" class to (1) and the values that correspond to "Negative" class to (0) as figure 9 shows.

```
# change values to numeric
data_df['classification'] = data_df['classification'].map({'Positive':1, 'Negative':0})

data_df.head(5)
```

Text	classification	clean text	text length	tokens
تلك الأيام " مغامرة سينمائية متميزة إنتظرت عرض فيلم تلك الأيام كثيرا ، خاصة لأنه مأخوذ عن رواية للاديب الكبير فتحي غانم و تحمل نفس الاسم الجذاب والتي تعتبر واحدة من أهم أعماله وعندما عرفت أن من سيقوم بإخراج الفيلم هو أحمد غانم ابن الروائي فتحي غانم شعرت بأنها شجاعة كبيرة منه خصوصا أن هذا العمل هو أولي تجاربه السينمائية وهو عمل ليس سهلا فتحويل عمل ادبي الى فيلم سينمائي يحتاج لخروج متمرس لذا كنت متخوف بعض الشيء من مستوى الإخراج ولكن وجود الفنان محمود حميدة في الفيلم جعلني أتوقع أنه سيكون فيلم جيد ولم يخيب حميدة توقعاتي فعند مشاهدته الفيلم أدركت أنني أمام عمل سينمائي متميز ومتخلف كثيرا بداية من مخرج خالف توقعاتي فقدم صوره متميزة للغاية ونجح في التعبير عن الحرب النفسية الموجودة بين الإطبال وانفسهم وذلك من خلال الاضواء وزوايا التصوير التي برع في استخدامها وكانت اروع المشاهد في الفيلم هي المظاهرات الخارجية التي يقوم بها احمد الفيشاوي فكانت علي اعلي مستوى سواء من ناحية التمثيل او التصوير او الإخراج. تدور قصة الفيلم حول شخصية سالم عبيد محمود حميدة وهو استاذ جامعي ومفكر سياسي وهو شخص من الممكن أن يفعل أي شيء ليروضي النظام وذلك طمعا في المنصب السياسي الكبير وهو أيضا نفس الشخص الذي تزوج فتاه لا تحبه وذلك طمعا في الحب الذي شاهدها وهو تعطيه إلى حبيبها السابق الذي توفي في حادث إرهابي وهو أيضا الشخص الذي ممن الممكن أن يفعل أي شيء للوصول إلى أهدافه . ويجب أن أشير أنه كان من		تلك الأيام " مغامرة سينمائية متميزة إنتظرت عرض فيلم تلك الأيام كثيرا ، خاصة لأنه مأخوذ عن رواية للاديب الكبير فتحي غانم و تحمل نفس الاسم الجذاب والتي تعتبر واحدة من أهم أعماله وعندما عرفت أن من سيقوم بإخراج الفيلم هو أحمد غانم ابن الروائي فتحي غانم شعرت بأنها شجاعة كبيرة منه خصوصا أن هذا العمل هو أولي تجاربه السينمائية وهو عمل ليس سهلا فتحويل عمل ادبي الى فيلم سينمائي يحتاج لخروج متمرس لذا كنت متخوف بعض الشيء من مستوى الإخراج ولكن وجود الفنان محمود حميدة في الفيلم جعلني أتوقع أنه سيكون فيلم جيد ولم يخيب حميدة توقعاتي فعند مشاهدته الفيلم أدركت أنني أمام عمل سينمائي متميز ومتخلف كثيرا بداية من مخرج خالف توقعاتي فقدم صوره متميزة للغاية ونجح في التعبير عن الحرب النفسية الموجودة بين الإطبال وانفسهم وذلك من خلال الاضواء وزوايا التصوير التي برع في استخدامها وكانت اروع المشاهد في الفيلم هي المظاهرات الخارجية التي يقوم بها احمد الفيشاوي فكانت علي اعلي مستوى سواء من ناحية التمثيل او التصوير او الإخراج. تدور قصة الفيلم حول شخصية سالم عبيد محمود حميدة وهو استاذ جامعي ومفكر سياسي وهو شخص من الممكن أن يفعل أي شيء ليروضي النظام وذلك طمعا في المنصب السياسي الكبير وهو أيضا نفس الشخص الذي تزوج فتاه لا تحبه وذلك طمعا في الحب الذي شاهدها وهو تعطيه إلى حبيبها السابق الذي توفي في حادث إرهابي وهو أيضا الشخص الذي ممن الممكن أن يفعل أي شيء للوصول إلى أهدافه . ويجب أن أشير أنه كان من		تلك الأيام " مغامرة سينمائية متميزة. إنتظرت عرض فيلم. تلك الأيام كثيرا. خاصة لأنه مأخوذ. عن رواية للاديب الكبير. فتحي غانم. و تحمل نفس الاسم الجذاب. والتي تعتبر واحدة من أهم أعماله. وعندما عرفت أن من سيقوم بإخراج الفيلم هو أحمد غانم. ابن الروائي فتحي غانم. شعرت بأنها شجاعة كبيرة منه خصوصا أن هذا العمل هو أولي تجاربه السينمائية وهو عمل ليس سهلا فتحويل عمل ادبي الى فيلم سينمائي يحتاج لخروج متمرس لذا كنت متخوف بعض الشيء من مستوى الإخراج ولكن وجود الفنان محمود حميدة في الفيلم جعلني أتوقع أنه سيكون فيلم جيد ولم يخيب حميدة توقعاتي فعند مشاهدته الفيلم أدركت أنني أمام عمل سينمائي متميز ومتخلف كثيرا بداية من مخرج خالف توقعاتي فقدم صوره متميزة للغاية ونجح في التعبير عن الحرب النفسية الموجودة بين الإطبال وانفسهم وذلك من خلال الاضواء وزوايا التصوير التي برع في استخدامها وكانت اروع المشاهد في الفيلم هي المظاهرات الخارجية التي يقوم بها احمد الفيشاوي فكانت علي اعلي مستوى سواء من ناحية التمثيل او التصوير او الإخراج. تدور قصة الفيلم حول شخصية سالم عبيد محمود حميدة وهو استاذ جامعي ومفكر سياسي وهو شخص من الممكن أن يفعل أي شيء ليروضي النظام وذلك طمعا في المنصب السياسي الكبير وهو أيضا نفس الشخص الذي تزوج فتاه لا تحبه وذلك طمعا في الحب الذي شاهدها وهو تعطيه إلى حبيبها السابق الذي توفي في حادث إرهابي وهو أيضا الشخص الذي ممن الممكن أن يفعل أي شيء للوصول إلى أهدافه . ويجب أن أشير أنه كان من

Figure 11. Convert value to numeric

- The next step was identifying the column that will contain the texts we want to predict their sentiment and the column that will contain the classification class label as figure 10 shows.

```
# idneitfy the data and the labels
data= data_df['clean text']
target= data_df['classification']
```

Figure 12. Identifying the column

- As figure 11 shows, we used TF-IDF vectorizer to extract the features from our data, and the vectorizer discovered 60467 words out of the 1368 rows of the given data.

```
1 # Use TfidfVectorizer for feature extraction (TFIDF to convert textual data to numeric form):
2 tf_vec = TfidfVectorizer()
3 X = tf_vec.fit_transform(data)
4 X.shape

(1368, 60467)
```

Figure 13. TF-IDF vectorizer

- 5- After that we plugged in our X array which contains the words after transforming them using TFIDF vectorizer, and our target array which contains the class labels, and we set up the test sample size to be 50% which leaves 50% for training as figure 12 shows.

```

1 # Training Phase
2 X_train, X_test, y_train, y_test = train_test_split(X, target, test_size=0.50, random_state=0)

1 print("Training: ", X_train.shape, y_train.shape)
2 print("Testing: ", X_test.shape, y_test.shape)

Training: (684, 60467) (684,)
Testing: (684, 60467) (684,)

```

Figure 14. Training phase

- 6- Finally, we plugged in our training and testing data to Multinomial Naïve Bayesian classifier function and got an accuracy of 70% as figure 13 shows.

```

1 # create the classifier and fit the training data and labels
2 classifier_nb = MultinomialNB().fit(X_train.todense(),y_train)
3
4 print("MultinomialNB accuracy: %.2f"%classifier_nb.score(X_test.todense(), y_test))
5 print('_'*100)
6
7 #do a 10 fold cross-validation
8 results_nb = cross_val_score(classifier_nb, X.todense(),target, cv=10)
9 print("\n10-fold cross-validation:")
10 print(results_nb)
11
12 print("The average accuracy of the MultinomialNB classifier is : %.2f" % np.mean(results_nb))
13 print('_'*100)
14
15 print("\nConfusion matrix of the MultinomialNB classifier:")
16 predicted_nb = classifier_nb.predict(X_test.todense())
17 print(confusion_matrix(y_test,predicted_nb))
18 print('_'*100)
19
20 print("\nClassification_report of MultinomialNB classifier:")
21 print(classification_report(y_test,predicted_nb))
22 print('_'*100)

```

MultinomialNB accuracy: 0.70

---

10-fold cross-validation:  
[0.69343066 0.69343066 0.69343066 0.69343066 0.69343066  
0.69343066 0.69343066 0.69852941 0.69852941]  
The average accuracy of the MultinomialNB classifier is : 0.69

---

Confusion matrix of the MultinomialNB classifier:  
[[ 0 202]  
[ 0 482]]

---

Classification\_report of MultinomialNB classifier:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	202
1	0.70	1.00	0.83	482
accuracy			0.70	684
macro avg	0.35	0.50	0.41	684
weighted avg	0.50	0.70	0.58	684

Figure 15. Classification report

As figure 15 shows, we plotted the True Positive Rate (TPR) and False Positive Rate (FPR) using the ROC Curve.

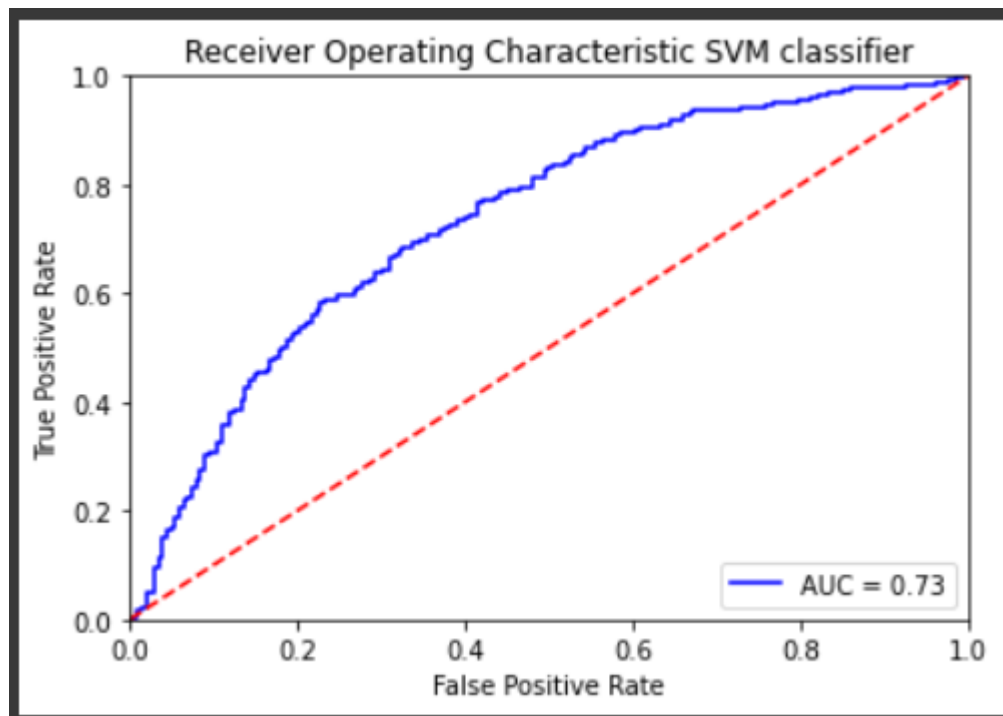


Figure 16. ROC Curve for our data

### **3. Discussion questions**

#### **3.1 Does the model appear valid and accurate on the test data?**

When looking at the results and comparing them with the evaluation techniques, we see that there is a discrepancy in the unbalanced data, so it does not seem valid, although the accuracy result is good.

#### **3.2 Does the model output/behavior make sense to the domain experts?**

We learned that it is not possible for the model to record a result of 100%, and although the model did not show any unexpected behavior or provide any unusual result, but the output was not very compelling.

#### **3.3 Do the parameter values make sense in the context of the domain?**

It can be expected that the model may not record very high results, especially in sentiment analysis data because it is rather complex, so we believe that there is a hidden problem regardless of the results and the scope on which it is.

#### **3.4 Is the model sufficiently accurate to meet the goals?**

Based on the outputs of our model, we see that the results are good and cover our goals, bearing in mind that we analyze people's feelings towards films, and that does not require very high accuracy, so the accuracy of our model is considered good, which is 70%.

#### **3.5 Are more data or inputs needed?**

Yes, when the data increases, the accuracy increases, and the model outputs will become better. In addition, we see that there is a discrepancy in the data, so a balance must be established between the data to see a much better result for our model.

## 4. Reference

[1] Rawat, A., 2022. *What is Descriptive Analysis?- Types and Advantages* | *Analytics Steps*. [online] Analyticssteps.com. Available at: <<https://www.analyticssteps.com/blogs/overview-descriptive-analysis>> [Accessed 8 August 2022].

[2] <https://www.javatpoint.com/machine-learning-naive-bayes-classifier>

[3] <https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>

[4] <https://numpy.org/doc/stable/user/whatisnumpy.html>

[5] [https://www.tutorialspoint.com/scikit\\_learn/scikit\\_learn\\_introduction.htm](https://www.tutorialspoint.com/scikit_learn/scikit_learn_introduction.htm)

[6] <https://www.activestate.com/resources/quick-reads/what-is-pandas-in-python-everything-you-need-to-know/>

[7] <https://www.activestate.com/resources/quick-reads/what-is-matplotlib-in-python-how-to-use-it-for-plotting/>