



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 1

#	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. Project Description.....	1
1.1 Introduction.....	1
1.2 Goal.....	1
1.3 Initial Hypothesis	1
1.4 Objective	2
1.5 Analysis Plan	2
2. Data Exploration	3
3. Data Issues	5
Issue#1: Punctuation.....	5
Issue#2: Text Normalization.....	6
Issue#3: English Characters.....	7
4. Tool Used.....	8
5. Reference.....	9

Table of figures

Figure 1.....	3
Figure 2.....	3
Figure 3.....	4
Figure 4.....	4
Figure 5.....	4
Figure 6.....	5
Figure 7.....	5
Figure 8.....	5
Figure 9.....	6
Figure 10.....	6
Figure 11.....	7
Figure 12.....	7

1. PROJECT DESCRIPTION

1.1 INTRODUCTION

In our world today, technology has become everywhere around us. It has more progress. The data grows quickly and has a large size, so storing, managing, and processing this data is very complex. As a result, we have the so-called 'huge data and cloud computing', indicating that huge data is a large amount of data of different types. Therefore, in this course, we learned about the concept of data analysis, which is the process of examining, purifying, converting, and modeling data to reveal useful information on a specific topic [1].

To obtain more experience in data analysis curricula, in our project, we aim to acquire knowledge about people's impressions and their satisfaction with the films they see by extracting and analyzing the available data on *Elcinema.com* on our topic and producing goals that we look forward to achieving.

1.2 GOAL

Nowadays, the industry market and the production of movies are developing very quickly. Films have become important in influencing society issues. *Elcinema.com*, a site that provides the largest Arabic cinematic content in which you can follow the latest news of the art world, and you can also book tickets online and watch movies and exclusive meetings with the stars of the Arab world where you can provide your evaluation and your observations for the movies that you have seen and read movie reviews for other site pioneers [2]. So, our main goal is to use our data analysis skills to explore and study the opinions of the site's pioneers to achieve the project goals.

1.3 INITIAL HYPOTHESIS

By the end of the project, we will either approve or disapprove the following hypothesis:

- People tend to watch movies of comedy nature.

If the hypothesis is approved, *Elcinema* will give an insight into the best comedy movies. On the contrary, if we cannot agree to the hypothesis, this may indicate that we must know what the site's pioneers want more.

- People tend to hate very long movies.

If the hypothesis is approved, the *Elcinema* will help to focus more on the number of movies and display movies that attract viewers. On the contrary, if we cannot agree to the hypothesis, this means that *Elcinema* should display films of all kinds and may not pay attention to the duration of the movie.

1.4 OBJECTIVE

Our objective is to understand how to implement data analysis techniques to produce a result that helps us know the behavior and interests of the viewers, and what they want to see more and what they like and what they do not like to give *Elcinema.com* some recommendations based on it, and at the end of this project, we hope that we will learn and apply data analysis techniques, starting From collecting data using Python libraries, to clean it, process it and even finally provide these data.

1.5 ANALYSIS PLAN

In our project we will implement the Data Analytics Lifecycle which defines the analytic process and best practices from discovery to project completion [3]. We will use "Google Collab" to write our code to clean the Dataset in addition to building the model and training it [9]. Data Analytics Lifecycle consists of 6 phases: Discovery, Data Preparation, Model Planning, Model Building, Communicate Results and Operationalize.

Step 1: Discovery

In this step, we got ready data on people's reviews from Elcinema.com.

Step 2: Data Preparation

In this step, we will clean the data by using preprocessing methods such as, ETL (Extraction, Transformation, Loading) and using Python libraries:

- Pandas: is a fast, powerful, flexible, and easy to use open-source data analysis and manipulation tool, built on top of the Python programming language [4].
- Re: A regular expression is a special sequence of characters that helps you match or find other strings or sets of strings, using a specialized syntax held in a pattern [5].
- String: Python String module contains some constants, utility function, and classes for string manipulation [6].
- Nltk: (Natural Language Toolkit) is an open-source Python library for Natural Language Processing [7].
 - Nltk tokenizer: Tokenizers divide strings into lists of substrings. we used tokenizers to splits a string using a regular expression, which matches either the tokens or the separators between tokens [8].

Step 3: Model Planning

In this step, we will define the methods and techniques that are based on our hypotheses and make sure that they meet the goals of our business by understanding data and relationships between variables.

Step 4: Model Building

In this step, we will implement the models specified in step 3, develop the dataset, and divide it into two parts:

1. Training Data: develop model in the Dataset
2. Test data: After completing the training, we will test the testing dataset to ensure that the model works well and as we want.

Step 5: Communicate Results

In this step, and after validating our building model we will communicate and document the key findings and determine if we succeeded or failed in our objectives.

2. Data Exploration

Data exploration is an important step to make pre-processing of data. Because it tells us a lot of information about the data we want to process and make us understand the nature of this data. Such as the data size, the types of information it contains and other important features that we can use to form appropriate cleaning processes.

- View the total number of data (rows) using len().

```
[21]
len(data_df)

1524
```

Figure 1. Exploration action (number of data)

- View the names of the columns included in the data to clarify the types of information included in the data using the columns.

```
[22] data_df.columns

Index(['Text', 'classification'], dtype='object')
```

Figure 2. Exploration action (Columns names)

- View some statistical data using describe().

```
[41] data_df.describe()
```

		Text	classification
count		1524	1519
unique		1522	3
top			Positive
freq			950

Figure 3. Exploration action (Statistical summary)

- View the count of non-null values contained in each column using count().

```
[42] data_df.count()
```

```
Text      1524
classification  1519
clean text  1524
tokens    1524
dtype: int64
```

Figure 4. Exploration action (columns non-null count)

- View each column with the corresponding non-null values count and its data type using info(). There is also a brief description of the datatypes included in the dataset.

```
[43] data_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1524 entries, 0 to 1523
Data columns (total 4 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
dtypes: object(4)
memory usage: 47.8+ KB
```

Figure 5. Exploration action (columns datatype and non-null values count)

3. Data Issues

Issue#1: Punctuation

Libraries used: pandas, and string

After examining the dataset, we noticed that the text has many unnecessary punctuations, so we decided to delete this punctuation. Figure 6 shows the list of punctuation that will be removed from the text in addition to the pre-specified punctuation that is provided with the string library using String.Punctuation constant [6].

```
arabic_punctuations = ' ' `÷×_—“”‘’!|+|~{ }',.?:"/, _][%^&*()_<>:;' '
english_punctuations = string.punctuation
punctuations_list = arabic_punctuations + english_punctuations
```

Figure 6. punctuation List

Figure 7 view the punctuation removal function that uses Re (Regular Expression) library to remove punctuation from a given text.

```
#this funcation take txet as input and return the text after removing punctuations
def remove_punctuations(text):
    translator = str.maketrans('', '', punctuations_list)
    return text.translate(translator)
```

Figure 7. remove punctuation function

Figure 8 view random sample before and after removing punctuation from the reviews.

data_df.sample(3)	
Text	classification
clean text	
<p>412 حلمي ما بين الدب باندأ والقرد ميسي نيج احمد حلمي مع مسلسل ميسي وتجسيد صوتي للقر ولأقي استحسن الكثير من الجمهور ولكنه فشل في فيلمه صنع في مصر بتجسيده لدميه دب باندأ حيث تصببه لعنه حينما تتمنى اخته الصغيره نور عثمان ان يتحول اخيها الى دمية وان يصبح الديدوب اخيها وبطريقه تقتقد الي اي حيكه دراميه يتم التحول ويحاول علاء احمد حلمي ان يعود الي شخصيته الادمية التي كانت شخصيه فاشله وينجح في النهاية ولكن بعد ان تتغير شخصيته ليصبح شاب يعتمد عليه اما بالنسبة للداء التمثيلي فتجد الطفلة نور عثمان التي لقت الانتظار في برنامج ارباب جود ايدل بادائها المتميز بالاستعراض والغناء لكن لم يتم الاستفاده من مواهبها وكذلك الحال مع الممثل المتألق في الفتره الاخيره بيومي فواد و النجمة (دلال عبد العزيز) فادوارهما هشة لا تليق بهما .. اما ياسمين رئيس ذات الوجه الملائكي والتي تألقت في (فتاة المصنع) فلم تات بجديد اما (انوار) فعلى الرغم من صغر الدور فانه جاء تغييراً له على الاقل كشكلاً</p>	<p>Positive</p>
<p>شمشون عنتر و ليلب عن فيلم شمشون عنتر و ليلب ياتري هل انت قوي الملاحظة و لا ينتفج و انت بتشرب حابه ساقعه و بتاكل فاكهه و يمكن كمان بتذاكر لعيالك لحد مايبيجي المدرس الخصوصوسي ما علينا المهم لو انت ركزت شويه هنلاقي كل من نطق اسم عنتر بالفيلم الصوت مش مناسب مع حركة الشايفاي دي اولاً ثانياً هنلاقي الصوت مش بتاع الممثل نفسه يعني لو عبد الوارث عسر المعلم رضوان يقول عنتر مش هنلاقي صوت عبد الوارث عسر</p>	

Figure 8. dataset before and after removing punctuation from the reviews

Issue#2: Text Normalization

Libraries used: pandas, and re

The Arabic text has many shapes for one letter. We will unify the appearance of these letters

Figure 9 view unify character's function that uses sub function in regex library to substitute the any occurrences of Arabic characters in the right with the one in the left in the given text [5].

```
def normalize_arabic(text):  
  
    text = re.sub("[اآإئإ]", "ا", text)  
    text = re.sub("ى", "ي", text)  
    text = re.sub("ه", "ه", text)  
    text = re.sub("ك", "ك", text)  
    text = re.sub("و", "و", text)  
  
    return text
```

Figure 9. Unify Characters Function

Figure 10 view random sample before and after removing punctuation from the reviews. For example, the word 'أجنبي' in the first line appeared after the text normalization app with the letter "ي" instead of the letter "ى", and the letter "ا" instead of the letter "إ".

data_df.sample(3)	
	Text classification
	clean text
1115	Neutral
<p>عندما يتحد الضعفاء عندما نشاهد أي فيلم أجنبي نجد طرفين طرف يمثل الخير تكمن قوته في الإرادة والعزيمة وطرف يمثل الشر يكمن في الخوف ودائما يكون الطرف الشرير هو الطرف الأقوي في البداية لأن الخوف من السهل جدا أن يستثار ويظهره أي شخص في العالم فلا يوجد انسان على وجه الأرض لا يخاف ولا يمتلك نقطة ضعف تكمن من هزيمته بسهولة بداية الفيلم تقليدية جدا شاهدنا مثلها من قبل في عدة أفلام مثلا فيلم الشيخ لأحمد عز وهي أن يستقيظ شخص فاقد الذاكرة في مكان لايعرف ماذا أحضره الى هنا وماذا حدث وماذا جرى له ؟ يكتشف بطلنا (جاك فروست) الأربعة الأقوياء حراس الكون وهي شخصيات خيالية معروفة لاي طفل في العالم (حورية الأستان و بابا نويل ورجل الأحلام وأخيرا أرنب عيد الفصح ولكل شخصية مكانها المعين وتمتلك قوتها الخاصة اضافة الي اتباعها المساعدين في الاعمال الموسية الي كل حارس وكما يوجد الخير يوجد الشر والذي يدعى بيتش بلاكأو الملعب الأسود ومن اسمه فهو يمثل الخوف ويمتلك هو الآخر جنوده من الفرسان الذي يمثلون الكوابيس لاي طفل ويسببون له الخوف وكلما زاد الخوف كلما ازدادت طاقة بيتش بلاك السلبية والتي تمكته من مواجهة الحراس يحاول (جاك الفروست) أن يكتشف لماذا تم اختياره كحارس من الحراس ويحاول سنتا كلوز (بابا نويل) أن يساعده أن يكتشف الطاقة الإيجابية التي بداخله ليتمكن من هزيمة خوفة هزيمة بيتش بلاك</p>	

Figure 10. Dataset before and after text normalization

Issue#3: English Characters

Libraries used: pandas, and re

One of the problems we have noticed in the dataset is that some reviews contain English characters, so to focus only on the Arabic text, we will remove all the English characters from the reviews, using the re (Regular expression) library in Figure 11, we remove any expression that matches the regular expression given.

```
def processPost(tweet):  
  
    #remove english letters  
    tweet= re.sub(r'[a-z]+'," ", tweet)  
    tweet= re.sub(r'[A-Z]+'," ", tweet)  
  
    return tweet
```

Figure 11. Remove English Characters Function

Figure 12 view a random sample before and after removing English characters from the reviews.

data_df["clean text"] =data_df["Text"].apply(lambda x: processPost(x))	
data_df.sample(3)	
Text classification	clean text
570	تغيير الماضي افلام العوده للماضي، وأمنية الرجوع للواء املا في تحقيق الأمنيات: من أكثر الأفلام التي يمكن تلاقي فيها نفسك، لأن ببساطة كل شخص فيها لو أتاحت له فرصة العوده للماضي وتغيير حياته، فيستغل الفرصة دي أفضل استغلال دي أفضل استغلال ممكن ويعتقدتش إننا هنفوتها كذا وفيلم من الأفلام التي بتتناول فكرة العوده للماضي أملا في تحقيق أمنية (Free Birds) ممكن ويعتقدتش إننا هنفوتها كذا. وفيلم ما. المخرج (جيمي هورار)، والمؤلفان (كريج مازن وديفيد اى ستورن) يقدمان توليفة بسيطة عن تلك التيمة من خلال قصة شيقه بتدور حول الديك الرومي (ريجى) الذي يحاول السفر للماضي وتحديدا عام 1621 مستخدما آلة الزمن حتى يتمكن من تغيير عادة البشر في أكل الديك الرومي ليلة عيد الميلاد..... ومن هنا فنلاحظ إن قصة الفيلم بتتناز بالبساطة، ويتركز من خلال حوار شيق على الكثير من النصائح والرسائل الهامة على المستوى العام سواء للأطفال أو الكبار في إطار كوميدي خفيف، فألتعاون ومحاولة الاتحاد لتحقيق غاية ومصلة جماعية تأتي في المقام الأول ثم المصلحة القومية. محاولة تسليان الخلافات الشخصية وعدم الالتفات لها أملا في تحقيق تلك الأمنية الجماعية. ويطعا من المعروف أن البطل الأول والأساسي في أي فيلم رسوم متحركة يعتمد على الجرافيك والألوان وتنفيذ المخاطر والمغامرات بحرفية وتقنية عالية، بالإضافة إلى الصوت والذي اعتقد أن أبطال العمل نجحوا بشكل جيد في توصيل تقاعلات وانفعالات الديك مع بعضها البعض ومع البشر، والمطاردات التي تعرضوا لها طوال أحداث الفيلم؛ فمثلا الممثل (أوين ويلسون) قدم دور الديك ريجي بطريقة مبهرة وقدر من خلال مستويات صوت مختلفة أنه يوصل سعي ريجي الدائم لإلقاء صنف الديوك من قاتمة الطعام في تلك الليلة، وتعاونه مع ريجي الدائم لإلقاء صنف الديوك من قاتمة الطعام في تلك الليلة وتعاونه مع صديقه الديك جاك وودي هاريلسون لتحقيق ذلك وقد حقق فيلم في أول أسبوع له بصالات عرض الولايات المتحدة الأمريكية حوالي مليون دولار أمريكي وقد ظل في المركز الثالث بالبيكس أوفيس الأمريكي لأكثر من أسبوعين متتاليين ثم بدأ في التراجع للخلف حتى المركز الخامس اعتقد أن فيلم المركز الخامس اعتقد أن فيلم من الأفلام التي ممكن تسبب أبك (Free Birds) متتاليين ثم بدأ في التراجع للخلف حتى المركز الخامس. اعتقد أن فيلم

Figure 12 Dataset before and after removing English Characters

4. Tool Used

These are the tools and libraries we have used during this phase.

1. Pands

A Python package providing fast, powerful, flexible, and easy to use open-source data analysis and manipulation tool, built on top of the Python programming language [4].

2. Re

A regular expression is a special sequence of characters that helps you match or find other strings or sets of strings, using a specialized syntax held in a pattern [5].

3. String

Python String module contains some constants, utility function, and classes for string manipulation [6].

4. NLTK

(Natural Language Toolkit) is an open-source Python library for Natural Language Processing [7].

5. Reference

- [1] “What Is Data Analysis? Methods, Techniques, Types & How-To,” *BI Blog |Data Visualization & Analytics Blog | datapine*, Mar. 09, 2022. <https://www.datapine.com/blog/data-analysis-methods-and-techniques/> (accessed August.04, 2022).
- [2] “What Is Elcinema Website” | *The Largest Arabic Movie Database | elcinema*, August. 04, 2022. <https://elcinema.com/> (accessed August.04, 2022).
- [3] W. Tian, Y. Zhao, *Optimized Cloud Resource Management and Scheduling*. Elsevier, 2014.
- [4] “pandas – Python Data Analysis Library,” *Pydata.org*, 2022. <https://pandas.pydata.org/> (accessed August.05,2022).
- [5] “re — Regular expression operations — Python 3.10.3 documentation,” *Python.org*, 2022. <https://docs.python.org/3/library/re.html> (accessed August. 05, 2022).
- [6] “string — Common string operations — Python 3.10.4 documentation,” *Python.org*, 2022. <https://docs.python.org/3/library/string.html> (accessed August. 06, 2022).
- [7] "NLTK - nltk package," nltk.org, 2022. [Online]. Available: <https://www.nltk.org/api/nltk.html> (accessed August. 06, 2022).
- [8] “tokenizers — splits a string into substrings using a regular expression— Python 3.10.4 documentation,” *Python.org*, 2022. <https://www.nltk.org/modules/nltk/tokenize/regexp.html> (accessed August. 06, 2022).
- [9] “Google colab” Colab.research, 2022. [Online]. <https://colab.research.google.com/> (accessed August. 05, 2022).