

Media Engineering and Technology Faculty  
German University in Cairo



# Arabic Sentiment Analysis System

Bachelor Thesis

Author: Ahmed Abdelghany Mohamed Ibrahim Madi  
Supervisors: Dr. Mohamed Ahmed El-Mahdy  
Submission Date: 15 May, 2017



Media Engineering and Technology Faculty  
German University in Cairo



# Arabic Sentiment Analysis System

Bachelor Thesis

Author: Ahmed Abdelghany Mohamed Ibrahim Madi  
Supervisors: Dr. Mohamed Ahmed El-Mahdy  
Submission Date: 15 May, 2017

This is to certify that:

- (i) the thesis comprises only my original work toward the Bachelor Degree
- (ii) due acknowledgement has been made in the text to all other material used

---

Ahmed Abdelghany Mohamed Ibrahim Madi  
15 May, 2017

# Acknowledgments

I begin by thanking Allah for giving me the strength and perseverance to complete this project till the very end, and for providing me with the knowledge needed to be able to work on it.

I would like to thank Dr.Mohamed El-Mahdy, for his encouragement and guidance throughout the period of this project,without which it would not have turned the way it did.

Furthermore, I must express my extreme gratitude to my parents and my family for providing me with continuous support and encouragement throughout the progress of my project. This work would not have been possible without them. Words are powerless to express my gratitude. Thank you.

# Abstract

Microblogging is one of the official communication platforms for internet users, thus becoming a rich domain for sentiment analysis and opinion mining. We focus on Arabic-speaking users due to the lack of focus on techniques devoted to Arabic sentiment analysis on data collected from Twitter, the most popular microblogging platform. This paper provides a comprehensive analysis on Arabic sentiment analysis with a data set of over 7,000 tweets classified into positive, negative and neutral by sentiment classifiers. The paper goes through the in depth analysis of methods used for data gathering, partitioning and linguistic analysis of the corpus and further explains the outcomes.

# Contents

<b>Acknowledgments</b>	<b>V</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Project Aim . . . . .	1
1.3 Thesis Organization . . . . .	2
<b>2 Background</b>	<b>3</b>
2.1 Tools and Technologies . . . . .	3
2.1.1 Sentiment analysis . . . . .	3
2.1.2 Twitter . . . . .	3
2.1.3 Arabic Sentiment Analysis . . . . .	5
2.1.4 Machine Learning . . . . .	5
2.1.5 Supervised Learning . . . . .	6
2.1.6 Python . . . . .	8
2.1.7 Scikit-Learn . . . . .	8
2.2 Previous Work . . . . .	8
2.3 Related Work . . . . .	9
<b>3 Implementation</b>	<b>11</b>
3.1 Dataset . . . . .	11
3.1.1 Dataset Collection . . . . .	11
3.1.2 Reading Dataset . . . . .	13
3.2 Feature Extraction . . . . .	14
3.2.1 Dataset Preparation . . . . .	14
3.2.2 Removing Stop-words . . . . .	14
3.3 Classifiers . . . . .	15
3.3.1 Naive Bayes . . . . .	15
3.3.2 Support Vector Machines . . . . .	17
3.3.3 Multi Layer Perceptron . . . . .	18
3.3.4 Random Forest . . . . .	18
3.3.5 Two-Step Classifiers . . . . .	18
3.3.6 Kollo ya waleed . . . . .	18

<b>4</b>	<b>Experimental Results</b>	<b>19</b>
<b>5</b>	<b>Conclusion</b>	<b>20</b>
<b>6</b>	<b>Future Work</b>	<b>21</b>
	<b>Appendix</b>	<b>22</b>
<b>A</b>	<b>Lists</b>	<b>23</b>
	List of Abbreviations . . . . .	23
	List of Figures . . . . .	24
	<b>References</b>	<b>27</b>



# Chapter 1

## Introduction

### 1.1 Motivation

Every day numerous users post profusely, which is one of the signs of their attachment to the platform. Starting with the definition of Microblogging, it is a combination of blogging and instant messaging that makes the communication process much easier compared to the old days when physical interactions had to be made or even traditional blogging. One of the reasons microblogging has managed to become popular is the less time spent developing content. This provides the user with the opportunity for frequent posts due to its character limit, which differentiates it from regular blogging. Also, it is the easiest way to share urgent or time-sensitive information. Further, microblogging is mobile convenient as it is not too hard to write a blog or post using a smartphone or tablet[26, 6]. As the platform becomes more popular, it enables higher number and frequency of posts. This high demand on microblogging allows the introduction of a new term, which is Sentiment Analysis.

### 1.2 Project Aim

To every project there is a main goal, and a subset of goals that branch from the target achieved. This project's main goal is to build a sentiment analysis system with the following criteria of having different features such as unigram, bigram and trigram as well as achieving the highest rate of accuracy classification possible. Meanwhile, the subsidiary goals are to tackle the obstacles created by pursuing sentiment analysis for Arabic that include the lack of availability of relevant content. Also, another subsidiary goal is to highlight the recent relevance of Arabic sentiment analysis with the growing usage of the Arabic language in microblogging platforms.

## 1.3 Thesis Organization

- Chapter 1- Introduction, in this chapter we discuss the motivation, project aim as well as an overview on how the thesis is organized.
- Chapter 2- Background, includes tools and technologies section that provides an overview about key concepts about Sentiment Analysis, Machine learning and the tools used to build this project like python and Scikit-Learn. Also, discusses some related and previous works related to this project.
- Chapter 3- Project Implementation gives a detailed description of the work done to build the project.
- Chapter 4- Shows the results of applying different approaches used in this project.
- Chapter 5- Conclusion
- Chapter 6- Future Work

# Chapter 2

## Background

### 2.1 Tools and Technologies

#### 2.1.1 Sentiment analysis

Sentiment analysis - otherwise known as opinion mining - is the process of discovering and determining the emotional meaning or tone behind a series of words. It can help in various ways. Firstly, by helping a company discover the public opinion towards their company or products. Said opinion aids in quality management, tactic and strategy planning as well as any marketing improvements whether through the business, economical or political position of the company and its products. All these changes are made based on the sentiment score provided by the sentiment analysis. Secondly, it can help political parties predict the public opinion towards them or the impact their campaign is having in order to align their goals with the public's interests. Moreover, the entertainment segment can benefit from sentiment analysis through collecting fans' feedback on the works of art by celebrities, authors and producers as well as opinions on their public interactions whether through interviews or social media[25].

When comparing the sentiment analysis to manual analysis or surveys, it should be mentioned that the human brain is the most accurate machine on earth. On the other hand, retrieving information in a systematic or computational way is more efficient than manual analysis and surveys as it is more time-efficient and cost-efficient with a good score that can be beneficial in receptive tasks[12].

#### 2.1.2 Twitter

There are multiple microblogging platforms including Jaiku and more recently Pownce. However, the most popular of microblogging platforms is Twitter. Twitter allows the

user to post statuses or updates that range to a maximum of 140 characters. By 2007, Twitter has managed to become one of the fastest rising microblogging platform with over 94,000 users. Twitter's popularity was due to the social communities formed through the mutual interest of its bloggers. The bloggers tended to stay for longer periods when they received positive comments and lasting social friendships through the platform. Also, the diversity of content provided whether through personal blogs or public blogs managed by celebrities in different entertainment segments that include movies, sports, music and art or politicians seeking public approval was enough to capture the attention of the public and engage in the microblogging phenomenon. Moreover, the fact that celebrities have verified accounts on twitter makes its information highly reliable as opposed to other gossip or events cyber venue[4]. 2.1 shows a verified account of the United States former president and a tweet from non-verified account of the US embassy about the Egyptian revolution.



**@USAbilAraby**  
USA bilAraby

#Egypt #Jan25 تعترف وزارة الخارجية الأمريكية بالدور التاريخي الذي يلعبه الإعلام الاجتماعي في العالم العربي ونرغب أن نكون جزءاً من محادثاتكم

8 Feb via web ☆ Favorite ↻ Retweet ↩ Reply

Retweeted by kabbani\_o and 26 others



(a) Non-verified account of United States embassy tweets about the Egyptian revolution .



(b) United States former president Barack Obama's Verified Account.

Figure 2.1: The diversity of content provided whether through personal blogs or public blogs managed by politicians.

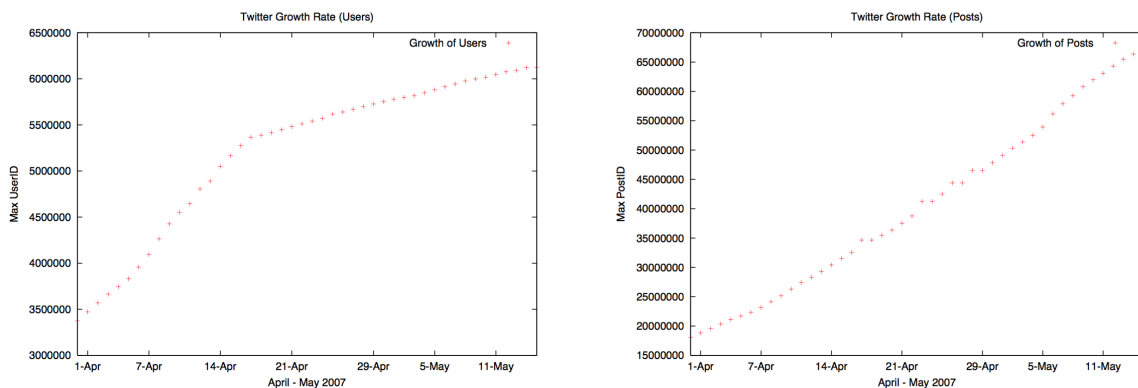


Figure 2.2: The number of Twitter users and posts during April-May,2007

Java et al. used Twitters API to show the rate of growth for users and posts for

two months in 2007. Although twitter was launched in 2006, it gains popularity after it won the South by SouthWest (SXWS) conference Web Awards6 in march,2007. Despite the number of new users has slowed, the number of posts has increased[3]. 2.2 shows the number of users and posts during April-May,2007 [3].

### 2.1.3 Arabic Sentiment Analysis

By 2016, Arabic was the fourth most commonly spoken language on Earth. The Arabic language is the official language for most Arabian countries including Egypt, Tunisia, Algeria, Saudi Arabia, Libya, Morocco and multiple gulf countries. It is widely spoken through the middle east and its relevance and importance is highlighted through the geographical and political affairs, as well as multiple works of art by authors such as Khalil Gobran that helped shape the Arab world in terms of culture. Also, the Islam's Holy Quran is in Arabic and with Muslims high percentage population of the world with 1/6th of the world's population, it creates a religious significance to the Arabic language.

Technologies have failed to incorporate the Arabic language in their work throughout the years, however more recently, many applications have used Arabic in search engines and document archiving tools, both known as core work. Furthermore, there has been a recent Arabic epidemic where the usage of the language has increased over the past few years. This increase is clear among social media platforms such as Facebook and for the sake of this topic, particularly Twitter. This further evolves the need for Arabic sentiment analysis and opinion mining[29].

### 2.1.4 Machine Learning

Machine learning enhances accuracy through allowing computers to modify or adapt their actions. These actions can include making predictions or controlling a robot. Said accuracy is evaluated and measured by how the correct actions are reflected in the chosen actions[22]. Machine learning recognizes human reasoning and repeats thought patterns and strategies that add to the decision-making process. It does so by producing classifying terms simple enough for humans to comprehend. Machine learning operations are resumed without the interference of humans, however background knowledge may be exploited to allow development of machine learning.[22]

A simple game of Scrabble can be used as an example. When you play against a computer, you start by winning. As the number of games progresses, the computer starts to win. The computer wins by understanding the pattern of strategies used by us and goes on to use these strategies against us and other players. By not starting from scratch every new player, it creates a form of generalization [22].

As mentioned, machine Learning aims to generate classifying expressions simple enough to be understood easily by the human. They must mimic human reasoning sufficiently to provide insight into the decision process. Like statistical approaches, background knowledge may be exploited in development, but operation is assumed without human intervention[11]. As show in 2.3 which describes the flow of a machine learning system .

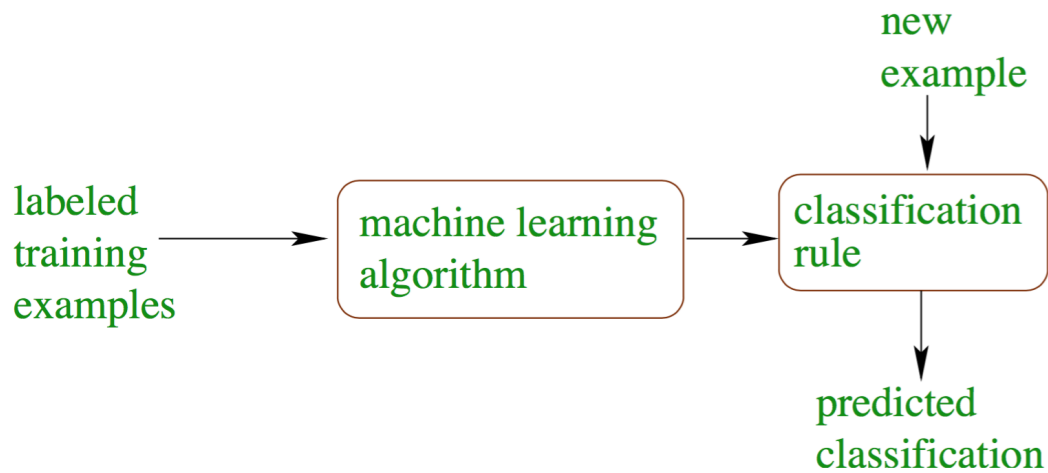


Figure 2.3: A diagram explaining the flow of machine learning problem.

### 2.1.5 Supervised Learning

There are two types of classification, one is known as Unsupervised Learning (or Clustering) and the other is Supervised Learning. Unsupervised learning is the classification of a data set through creating classes or clusters in the data. Meanwhile, Supervised Learning already acknowledges the presence of classes or clusters and aims to establish a rule that classifies newly introduced observations into one of the existing classes. The term "discrimination" is also used when talking about Supervised Learning due to correctly constructing a classification rule from formerly classified data. If former data is classified with a high rate of accuracy, the assumption is that someone (The Supervisor) is able to classify without error[31].

The function of supervised learning is derived from the training data, which is assembled from training examples. Every example is a pair of input and output, where input is usually a vector and output is the desired result, also known as the supervisory signal. According to (citation article supervised learning) the inferred function is called a classifier. The classifier should be able to generalize the training data in order to for

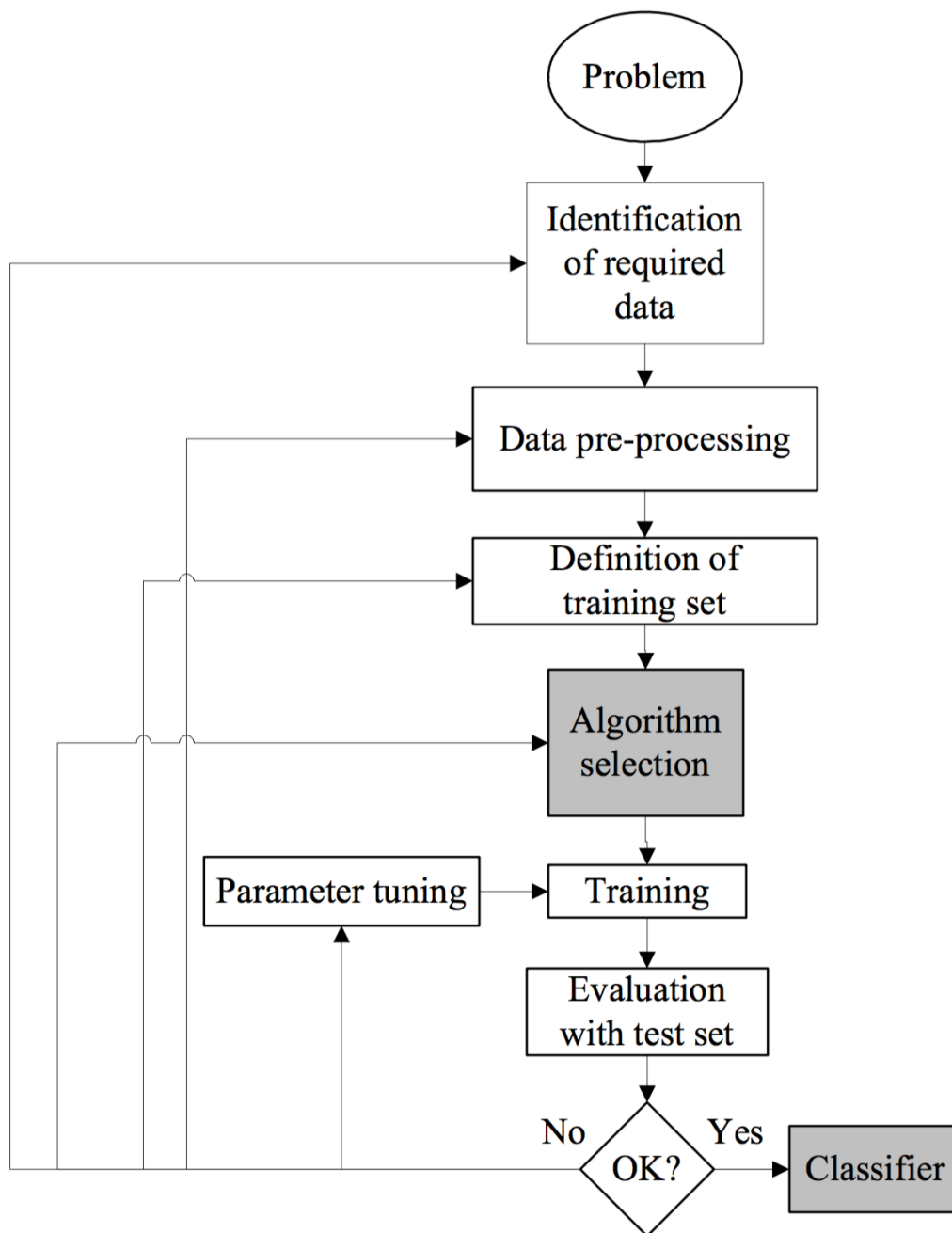


Figure 2.4: The supervised Learning Process

an unknown input to have a correctly predicted output[11] as shown in 2.4[17].

### 2.1.6 Python

According to Dubois Milmann and Avaizis, One of the most common languages used nowadays in scientific computing is Python[23, 13]. This goes back to its high-level interactive nature and its maturing ecosystem of scientific libraries. It is considered to be the first choice for algorithmic development and exploratory data analysis. Moreover, as a general-purpose language, it is increasingly used not only in academic settings but also in industry settings[15].

### 2.1.7 Scikit-Learn

Scikit learn is a python toolbox that manages to provide state-of-the-art implementations of multiple popular machine learning algorithms. All the while, maintaining a user-friendly interface in line with the Python language. The need for statistical data analysis is satisfied due to the friendly nature of the interface that allows many stakeholders to use it. The stakeholders include not only non-specialists in the software and web industries but also fields outside of computer science, such as Biology or Physics[15]. There are many features that differentiate Scikit learn from other machine learning toolboxes in Python, some of which include, firstly, the BSD license under which it is distributed. Secondly, when compared to MDP and Pybrain, Scikit learn incorporates compiled code for better efficiency (Zito et al., 2008)[32]. And while Pymvpa has optional dependencies such as R and shogun, and iv, Scikit learn only depends on numpy and scipy to facilitate easy distribution (Hanke et al., 2009)[19].

Also, another comparison of Pybrain to Scikit learn is the fact that Pybrain uses a data-flow frame work while Scikit learn focuses on imperative programming. While the package is mostly written in Python, it incorporates the C++ libraries LibSVM (Chang and Lin, 2001) and LibLinear (Fan et al., 2008) that provide reference implementations of SVMs and generalized linear models with compatible licenses. Binary packages are available on a rich set of platforms including Windows and any POSIX platforms[9, 27]. Furthermore, thanks to its liberal license, it has been widely distributed as part of major free software distributions such as Ubuntu, Debian, Mandriva, NetBSD and Macports and in commercial distributions such as the Enthought Python Distribution[15].

## 2.2 Previous Work

As was mentioned in 2.1.1, according to ( Pang and Lee the increase of blog and social network usage is directly proportional to the interest of researchers in opinion mining and sentiment analysis. They continue to discuss several techniques and strategy approaches for opinion-oriented information collection[25]. Although other works have ignored microblogging when discussing sentiment analysis, Yang et al used web-blogs data where



the users mood was determined through emotion icons assigned to blog posts and related comments. Next, for their research, sentiments were classified at sentence level through SVM and CRF learners. Various strategies were then discussed and analysed to set the overall sentiment of the document, when said sentiment matches the sentiment of the last sentence of the document, the strategy used is announced as the chosen strategy[10]

Another study by Read, 2005 stated that emoticon included texts were collected with the purpose of forming a training set for sentiment classification. The data was collected from Usenet newsgroups. The dataset that was divided into positive subset and negative subset where the positive sample included texts with happy emoticons, while the negative samples included sad or angry emoticons. A 70% accuracy rate was achieved when SVM and Nave Bayes, both known as Emoticons-trained classifiers, were used.[28]

Go et al., 2009 used the same approach , similar to the one performed by Read 2005, where the data set was classified similarly with an 81% rate of accuracy through the Nave Bayes classifier. On the other hand, it showed poor performance results when a third class known as neutral was introduced[5].

(Schmid ,1994) states that another method whose main interest is a difference of tags distribution between sets of text in English. It starts by checking the distribution of word frequency. Next, TreeTagger is used for English to tag all posts. This is a highlighted feature of the method due to the particular interest in difference of tags distribution between texts (positive, negative and neutral). Then, in order to train the sentiment classifier, features are extracted from the dataset acquired[30]. Nevertheless, (Pang et al 2002) have achieved better results when they used term presence instead of term frequency[8].

Whilst carrying out classification of movie reviews, Pang et al, 2002 have stated that uni-grams have greater success rate than bigrams[8]. Meanwhile, Dave et al. argue that bigrams and trigrams work better for the product review polarity classification[18]. Others recently working on sentiment analysis for Arabic include (Nabil, Aly, Atiya, 2015). With over 84,000 tweets collected, leaving 54,176 Arabic tweets after filtration, they have divided their dataset into 4 data sets with positive, negative, subjective and subjective mixed categories. All the while, performing a standard partitioning to the data set using a wide range of standard classifiers to perform 4 way sentiment classification to avoid sentiment polarity classification problems[21].

## 2.3 Related Work

Although the detection of user sentiment in text is recent and even harder to find in Arabic, there are still some articles that discuss sentiment analysis for Arabic. There are various processes for data collection and performing sentiment analysis, some of which include the SAMAR system, where multi domain datasets are exposed to subjectivity and sentiment analysis for Arabic social media. Said datasets are collected from Wikipedia TalkPages, Twitter and Arabic Forums. (Abdul Mageed et al., 2014)[2]. Similar to

Nabil, Aly, Attiya, 2013, (El Sahar and El Beltagy, 2015) have also classified data into four datasets. This was intended for the purpose of building a multidomain Arabic resource also known as sentiment lexicon[14]. All the while, (Nabil et al., 2014) and (El Sahar and El Baltagy, 2015) introduced a semi-supervised approach for building a sentiment lexicon that when used appropriately can be effective and efficient in sentiment analysis[14, 20].

Another article by (Aly and Atiya, 2013) presented LABR for a book reviews dataset collected from GoodReads[7]. Another article dataset related to entertainment was discussed by (Rushdi-Saleh et al., 2011) where 500 movie reviews were collected from multiple webpages and used as datasets for sentiment analysis. Moving on to works relevant to microblogging, (Refaee and Rieser, 2014) have managed to collect a dataset of 8,868 tweets that was manually annotated Arabic social corpus and they continue to review their method of collecting and annotating the corpus.

# Chapter 3

## Implementation

### 3.1 Dataset

#### 3.1.1 Dataset Collection

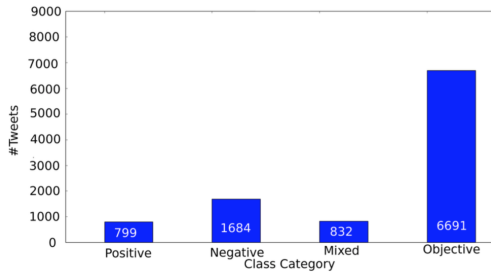
##### ASTD

Arabic Sentiment Analysis Tweets Dataset (ASTD) was collected by Nabil et al. have used. They divided their work to collect those tweets into three stages. Also, they have used high-tech tools to collect those tweets. Starting with the first stage, which aimed to determine the most active Egyptian Twitter accounts with their recent tweets till November 2013. They ended up with almost 30 accounts and about 36,000 tweets[21]. Proceeding to the second stage, they have crawled EgyptTrends 3, which is a Twitter page for the top trending hash-tags in Egypt. They have got a list of 2500 unique hash-tags which have used again to download more tweets. In this process, the number of tweets has increased to reach 84,000 tweets. Finally, they filtered those tweets and removed the non-Arabic tweets and redundant expressions like HTML, and ended up with 54,716 Arabic tweets[21]. 3.1 shows some examples of the collected dataset.[21]

After collecting those tweets, they moved on to another process which is annotating those tweets. Annotation is the process of adding information to already existing data. Whenever there is a large number of semi or unstructured data, Annotation is the best method to use to convert those unstructured data into a well-structured data [21]. Also, another high-tech is used to automate this process. Amazon Mechanical Turk (AMT) service provides an API called Boto 4 used to annotate those tweets. They have used four categories through this operation, which are: Objective, Subjective Positive, Subjective Negative and Subjective mixed. When a tweet was assigned to three or more labels, it was considered useless and was ignored. Another situation that has been accepted for further processing and was considered conflict-free was occurred when a tweet was

	<b>Tweet</b>	<b>Translation</b>	<b>Rate</b>
1	أكثر شعور بوجع ! #لما تجوع في بيت مو بيتكم **	Feeling that hurts ^ ! #To starve in a house not yours	Negative
2	محبين البرنامج يزدوا :)	Fans of El-Bernameg are increasing :)	Positive
3	#كفاية اسفاف	#stop smallness	Negative
4	الطاقة البشرية اذا ما احسن استغلالها هي رصيда و ليست عبئا قوتنا في عددنا	Human energy if properly exploited is an asset and not a burden our strength in our numbers	Positive
5	احبي الشيخ حسن عبد البصير امام مسجد سيدى جابر الذي رفض تعليمات الأوقاف بنفاق مرسى في خطبة الجمعة تعلموا الاستقامة أبها #الاخوان الكاذبون	I greet Sheikh Hassan AbdelBassir Imam Sidi Gaber mosque, who refused the instructions of the endowments to hypocrite Morsi in his Friday sermon learn the integrity liars brotherhood	Mixed
6	هل يُنوح ايتكو مدريد بلقب الليجا الأحد القادم؟ #برشلونة	Is Atletico Madrid going to be crowned La Liga next Sunday? # Barcelona	Objective

Figure 3.1: ASTD Tweets Examples



(a) The number of Tweets for each category.

Total Number of conflict free tweets	10,006
Subjective positive tweets	799
Subjective negative tweets	1,684
Subjective mixed tweets	832
Objective tweets	6,691

(b) Tweets Histogram.

Figure 3.2: ASTD Statistics.

assigned to two labels. Finally, they have ended up with 10,000 labeled tweets[21]. 3.2 shows the number of collected Tweets per each class[21].

### The Tweets dataset

Absualla et al. worked on the same topic and used a tweet crawler. Unfortunately, they did not mention it. They have collected 2,000 labeled tweets (1,000 positive and 1,000 negative ones) which were written in both Modern Standard Arabic (MSA) and the Jordanian dialect. They did not stick to one topic, they tried to cover various topics such as politics, arts and sports. Also, they did not use any automated system to annotate the data. Tweets were annotated using human experts. Exactly two human experts labeled those tweets. During the process of categorizing them, if they agreed on the same label of a tweet, it was assigned directly to the respective class. If a conflict has occurred, a third expert was used to label this tweet[24].

After collecting the data, they moved to another procedure which is dataset preparation. They have started with separating each tweet in a single file and split those tweets to be in a positive or negative folder. Then, they have used a naive algorithm to remove

Positive	انا احب هذا الكاتب I like this author
Negative	الله يكون في عون الفقير و الطبقة المتوسطة سوف تنحدر أكثر و أكثر God help the poor and the middle class will diminish more and more

(a) The Tweets Dataset Examples.

	Positive	Negative
Total tweets	1000	1000
Total words	7189	9769
Avg. words in each tweet	7.19	9.97
Avg. characters in each tweet	40.04	59.02

(b) The Tweets Dataset Statistics

Figure 3.3: The Tweets Dataset Examples and statistics

the repeated letter (e.g., muuuuuch = ) and the MS Word dictionary as a reference to correct misspelled words. The naive algorithm worked as follow; it started with counting the number of letters in each word and checked whether this number exceeded five. If a letter's occurrences exceeded five, it eliminated the repeated letters and looked it up in the MS word dictionary[24]. 3.3 shows examples of the collected dataset and some statistics about it.[24]

Finally, they tried to enhance the dataset's quality by removing stop-words. A combined list of stop-words, which was obtained from the Khoja stemmer tool and added manually from different Arabic dialects was used to remove stop-words from the whole tweets. Another enhancement was used, which normalized some letters to one shape.

### 3.1.2 Reading Dataset

Those were the datasets that I worked on. Now, moving on to the first step of my implementation, which was reading those datasets correctly and modifying them.

#### ASTD

ASTD was represented as a file containing the whole tweets with its labels. Each tweet is represented by one line with its corresponding label. The first part of each line is the tweet itself and the rest is the label of this tweet. The aim of this part was to read the tweets and their labels without any redundant expressions (e.g., the new line and tab expressions) and combine the objective and subjective mixed classes into one "Neutral" class.

#### The Tweets Dataset

The Tweets dataset was represented as folder containing two files, the first one is positive which contains the positive tweets and the another one is negative which contains the negative tweets as it was mentioned in 3.1.1. The procedure went the same way as ASTD one, except that there are only two classes. There is no need to combine two classes into one class. Yet, there were some bad encoding of some tweets which was ignored during the reading process.

## 3.2 Feature Extraction

Feature extraction is considered to be the most difficult part in sentiment analysis as it is the major procedure that every procedure depends on. Feature extraction is the process of transforming the input data into a set of features. Features are a set of distinctive properties of the input that helps in differentiating between the categories of input patterns. Sentiment classification over Twitter is usually affected by the noisy nature (stop-words, punctuation and hashtags) of tweets data. In this section, different approaches of feature extraction will be discussed. In addition to some useful datasets' preparation techniques.

### 3.2.1 Dataset Preparation

#### Removing Hash-tags

Hash-tags is a major issue faced when dealing with Tweets' datasets. The hash-tag symbol is commonly used by people before a relevant keyword or phrase in their Tweet to categorize those Tweets or to search for related topics or Tweets to this hash-tag. Hash-tagged words that become very popular are often Trending Topics. For the sake of our topic Sentiment Analysis, Hash-tags are redundant expressions that could be useless or might decrease the accuracy of the models as it can be included in most of those Tweets. So, removing them is considered to be a good start for feature extraction. This work is done by the "remove-hash-tags" function which takes all the words in those Tweets as an input, then it proceeds with removing any hash-tag symbol that comes before any word.

#### Removing Punctuation

There is no doubt that punctuation is used to create sense, clarity and stress in sentences. Punctuation marks are used to structure and organize what you want to write. Those signs and symbols are given to the reader to show how a sentence is constructed, how it should be read and make the meaning clear. On the other hand, sentiment analysis does not take into account the effect of those marks. It may weaken the system and without those marks, the system performance will be more efficient. The remove-punctuation function handles those marks, which convey no meaningful sentiment and removes them wherever they are found.

### 3.2.2 Removing Stop-words

A well-known method that helps to reduce the feature space of the classifiers and enhance the efficiency of the system to give more accurate results, is the removal of stop-words. This process aims to discard non-discriminative words which have no sentimental meaning. (Saif et al, 2104) showed that applying those methods to Twitter in the context of

sentiment analysis obtaining contradictory results. Regarding this issue, the effectiveness of removing stop-words in the context of Twitter sentiment classification has been debated in the last few years.. (Bakliwal et al., 2012; Pak and Paroubek, 2010; Zhang et al., 2012; Speriosu et al., 2011; Gokulakrishnan et al., 2012; Kouloumpis et al., 2011; Asiaee T et al., 2012) support the idea of removing those stopwords which have no sentimental meaning. On the other hand, Saif et al., 2012b; Hu et al., 2013b; Martnez-Camara et al., 2013; Hu et al., 2013a) believes that removing those stop-words weaken the system and decrease the efficacy of the system as those words indeed carry sentiment information[16].

### 3.3 Classifiers

#### 3.3.1 Naive Bayes

In the supervised learning algorithms known as nave Bayes methods, the Bayes theorem is applied while every pair of features is naively assumed to be independent. As such, the theorem states:

$$P(y|x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n|y)}{P(x_1, \dots, x_n)} \quad (3.1)$$

Where  $y$  is a class variable and  $x_1$  through  $x_n$  is the dependent feature vectors. Since the nave assumption is mathematically represented as shown in 3.2, the relationship could be boiled down to 3.3 for all  $i$ .

$$P(x_i|y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i|y) \quad (3.2)$$

$$P(y|y, x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i|y)}{P(x_1, \dots, x_n)} \quad (3.3)$$

Given the input,  $P(x_1, \dots, x_n)$  is a constant. Subsequently, this classification rule could be applied:

$$P(y|y, x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i, y) \quad (3.4)$$

↓

$$\hat{y} = \operatorname{argmax} P(y) \prod_{i=1}^n P(x_i, y) \quad (3.5)$$

Following this logic, Maximum A posteriori (MAP) could be employed to estimate the relative frequency of the class  $y$  in a given training set,  $P(y)$ , as well as  $P(x_i|y)$  [1].

### Bernoulli

As mentioned, different Nave Bayes (NB) methods have different assumptions regarding the distribution of  $P(x_i|y)$ . In this regard, BernouliNB considers the data to follow a multi-variate Bernoulli distribution. Regardless of the number of features, each is a binary variable. Any other kind of data is made into a boolean variable using the binarize parameter. The decision rule in this case is 3.6. In case the input data is in the form of text, word frequency, instead of word count, could be employed when this classifier is being trained and, subsequently, used[1].

$$P(x_i|y) = P(i|y)x_i + (1 - p(i|y))(1 - x_i) \quad (3.6)$$

### Multinomial

In case the data is multinomially distributed, MultinomialNB could be implemented. It is classically used with text input data represented by word count vectors. Unlike BernoulliNB, which penalizes a non-occurring feature, MultinomialNB merely ignores such features. To parameterize the distribution, vectors  $\theta_y = (\theta_{y1}, \dots, \theta_{yn})$  are considered for each class  $y$ , where  $n$  is the number of considered features (for example, the size of the vocabulary in text-input data) and  $\theta_{yi}$  is the probability  $P(x_i|y)$  of having a feature  $i$  in a sample of class  $y$ . In this context, a smoothened version of maximum likelihood, like counting relative frequencies, could be used to estimate the parameters  $\theta_y$  as shown in 3.7. Where  $N_y$  is the count of all features in a class  $y$  and  $N_{yi}$  is the frequency of a given feature  $i$  in a sample belonging to class  $y$  in the training dataset  $T$  [1].

$$\hat{\theta}_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha n} \quad (3.7)$$

As mentioned, MultinomialNB ignores non-occurring features and this, the smoothing priors  $\alpha \geq 0$  removes zero probabilities from further computations. Laplace smoothing sets  $\alpha$  to 1 whereas Lidstone smoothing takes  $\alpha$  to be  $> 1$  [1].

### Gaussian

Finally, GuassianNB is suited to classifications where the feature likelihood  $P(x_i|y)$  follows a Gaussian distribution. Under this assumption 3.8. Where maximum likelihood is used to give an estimate to each of the parameters  $\sigma_y$  and  $\mu_y$  [1] .

$$P(x_i, y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right) \quad (3.8)$$



### 3.3.2 Support Vector Machines

Support Vector Machines (SVM) is one of the most popular algorithms in modern machine learning. It was introduced by Vapnik in 1992 and has taken radically since then, principally because it often provides very impressive classification performance on reasonably sized datasets. SVM approach based on finding the optimal separating hyperplane between two classes by maximizing the margin between them. Points which laying on the boundaries considered support vectors and middle of the margin is the optimal separating hyperplane. According to Thorsten Joachim SVM approach is one of the best approaches to be used in sentiment analysis, because it can deal with high dimensional input space. Also, SVMs use over-fitting protection which does not depend on the number of the input features. Finally, most of the text classification problems are linearly separable[22].

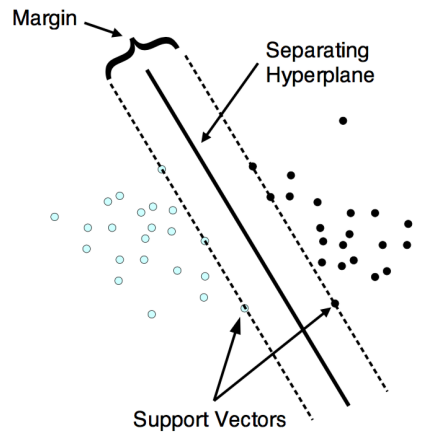


Figure 3.4: SVM

#### SVM Multi-Class Classification

A SVM classifier which considers a binary classifier, that is, the class labels can only be just two values. In the real world, many problems have more than two classes. However, many approaches tried to solve this problem. One of them is the SVM M-class classification, which based on building M-classifiers. Each classifier works on separating one class from the M-1 classifier, which known as "One Versus Rest" classification. Finally, before applying the function of SVM combine those M classifiers.

### **3.3.3 Multi Layer Perceptron**

### **3.3.4 Random Forest**

### **3.3.5 Two-Step Classifiers**

### **3.3.6 Kollo ya waleed**

## Chapter 4

# Experimental Results

# Chapter 5

## Conclusion

Conclusion

# Chapter 6

## Future Work

Text

# Appendix

# Appendix A

## Lists

# List of Figures

2.1	The diversity of content provided whether through personal blogs or public blogs managed by politicians. . . . .	4
2.2	The number of Twitter users and posts during April-May,2007 . . . . .	4
2.3	A diagram explaining the flow of machine learning problem. . . . .	6
2.4	The supervised Learning Process . . . . .	7
3.1	ASTD Tweets Examples . . . . .	12
3.2	ASTD Statistics. . . . .	12
3.3	The Tweets Dataset Examples and statistics . . . . .	13
3.4	SVM . . . . .	17



# Bibliography

- [1] Naive bayes. [http://thhttp://scikit-learn.org/stable/modules/naive\\_bayes.html#gaussian-naive-bayes](http://thhttp://scikit-learn.org/stable/modules/naive_bayes.html#gaussian-naive-bayes).
- [2] Muhammad Abdul-Mageed and Mona Diab. AWATIF: A Multi-Genre Corpus for Modern Standard Arabic Subjectivity and Sentiment Analysis. *In LREC.*, pages 3907–3914, 2012.
- [3] Tim Finin Akshay Java, Xiaodan Song and Belle Tseng. Why We Twitter: Understanding Microblogging Usage and Communities. *Joint 9th WEBKDD and 1st SNA-KDD Workshop 07, San Jose, California , USA*, 2007.
- [4] Tim Finin Belle Tseng Akshay Java, Xiaodan Song. Why We Twitter: Understanding Microblogging Usage and Communities. *Proceedings of the Ninth WebKDD and First SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis (WebKDD/SNA-KDD '07) (pp. 5665). New York: ACM Press.*
- [5] Lei Huang Alec Go and Richa Bhayani. Twitter sentiment analysis. *Final Projects from CS224N for Spring 2008/2009 at The Stanford Natural Language Processing Group.*, 2009.
- [6] Patrick Paroubek Alexander Pak. Twitter as a Corpus for Sentiment Analysis and Opinion Mining. *Universite de Paris-Sud, Laboratoire LIMSI-CNRS, Laboratoire LIMSI-CNRS, Orsay Cedex, France.*
- [7] Mohammed Aly and Amir Atiya. abr: Large scale arabic book reviews dataset. *In Meetings of the Association for Computational Linguistics (ACL), Sofia, Bulgaria.*, 2013.
- [8] Lillian Lee Bo Pang and Shivakumar Vaithyanathan. sentiment classification using machine learning techniques. *In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP).*, pages 79–86, 2002.
- [9] C.C. Chang and C.J. Lin. LIBSVM: a library for support vector machines. <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 2001.

- [10] Kevin Hsin-Yih Lin Changhua Yang and Hsin-Hsi Chen. Emotion classification using web blog corpora. In *WI 07: Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence. Washington, DC, USA. IEEE Computer Society.*, pages 275–278, 2007.
- [11] C.C. Taylor D. Michie, D.J. Spiegelhalter. *Machine Learning, Neural and Statistical Classification*. 1994.
- [12] Kalina Bontcheva Diana Maynard. Challenges of Evaluating Sentiment Analysis Tools on Social Media. *University of Sheffield, Department of Computer Science Regent Court, 211 Portobello, Sheffield, S1 4DP, UK*.
- [13] P.F. Dubois. Python: Batteries Included. *IEEE/AIP*, 9, 2007.
- [14] Hady ElSahar and Samhaa R El-Beltagy. Building large arabic multi-domain resources for sentiment analysis. In *Computational Linguistics and Intelligent Text Processing.*, pages 23–34, 2015.
- [15] Alexandre Gramfort Vincent Michel Bertrand Thirion Olivier Grisel Mathieu Blondel Peter Prettenhofer Ron Weiss Vincent Dubourg Jake Vanderplas Alexandre Passos David Cournapeau Matthieu Brucher Matthieu Perrot Fabian Pedregosa, Gael Varoquaux and Edouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 2012.
- [16] Yulan He Hassan Saif, Miriam Fernandez and Harith Alani1. On Stopwords, Filtering and Data Sparsity for Sentiment Analysis of Twitter. In: *Proc. 9th Language Resources and Evaluation Conference (LREC), Reykjavik, Iceland.*, 2014.
- [17] S.B. Kotsiantis. Building large arabic multi-domain resources for sentiment analysis. *Informatica 2007.*, 31, 2007.
- [18] Steve Lawrence Kushal Dave and David M. Pennock. Mining the peanut gallery: opinion extraction and semantic classification of product reviews. In *WWW03: Proceedings of the 12th international conference on World Wide Web.*, pages 519–581, 2003.
- [19] P.B. Sederberg S.J. Hanson J.V. Haxby M. Hanke, Y.O. Halchenko and S. Pollmann. PyMVPA: A Python toolbox for multivariate pattern analysis of fMRI data. *Neuroinformatics*, 7(1):37–53, 2009.
- [20] Mohamed A. Aly Mahmoud Nabil and Amir F. Atiya. Supervised machine learning: A review of classification techniques. *LABR: A large scale arabic book reviews dataset.*, 2014.
- [21] Mohamed Aly Mahmoud Nabil and Amir F. Atiya. ASTD: Arabic Sentiment Tweets Dataset. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal*, pages 2515–2519, 2015.

- [22] Stephen Marsland. *MACHINE LEARNING An Algorithmic Perspective*.
- [23] K.J. Milmann and M. Avaizis. Scientific Python. *IEEE/AIP*, 11, 2011.
- [24] Mohammed A. Shehab Nawaf A. Abdulla, Nizar A. Ahmed and Mahmoud Al-Ayyoub. Arabic Sentiment Analysis: Lexicon-based and Corpus-based. *IEEE Jordan Conference on Applied Electrical Engineering and Computing Technologies (AEECT)*., 2013.
- [25] Bo Pang and Lillian Lee. Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(12):1–135, 2008.
- [26] Apoorv Agarwal Boyi Xie Ilia Vovsha Owen Rambow Rebecca Passonneau. Sentiment Analysis of Twitter Data. *Department of Computer Science, Columbia University, New York, USA*.
- [27] C.J. Hsieh X.R. Wang R.E. Fan, K.W. Chang and C.J. Lin. LIBLINEAR: a library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [28] Jonathon Read. Using emoticons to reduce dependency in machine learning techniques for sentiment classification. *In ACL. The Association for Computer Linguistics*, 2005.
- [29] Ahmed Ali Samhaa R. El-Beltagy. Open issues in the sentiment analysis of Arabic social media: A case study. *IEEE Innovations in Information Technology (IIT), 2013 9th International Conference*, 2013.
- [30] Helmut Schmid. Probabilistic part-of-speech tagging using decision trees. *In Proceedings of the International Conference on New Methods in Language Processing*., pages 44–49, 1994.
- [31] Reddi Sanjeev Kumar Shirdi Wazeed Baba<sup>1</sup>. Data Mining: Text Classification System for Classifying Abstracts of Research Papers. *International Journal of advance Research, Ideas and Innovations In Technology*, 2, 2016.
- [32] L. Wiskott T. Zito, N. Wilbert and P. Berkes. Modular toolkit for data processing (MDP): A Python data processing framework. *Frontiers in Neuroinformatics*, 2, 2008.