

Requirements Specification Document

Project Title: Sentiment Analysis and Opinion Mining of the Arabic Web (Digital Content)

Selected ITAC Program: Advanced Research Project (ARP)

Milestone: 3

Academic and ICT Industry Partners

Organization Name	Contact Name	Role
American University in Cairo	Ahmed Rafea	Professor at CSE Department
LINK Development	Hanan Abdel Meguid	Chief Executive Officer

AUC Research Team

Name	Contact Details	Role	Date
Prof. Ahmed Rafea	rafea@aucegypt.edu	Principal Investigator	3/2011-present
Nada Ayman	nadaaym@aucegypt.edu	Researcher A	3/2011-7/2011
Islam Elnabarawy	islam.o@aucegypt.edu	Researcher A	3/2011-10/2011
May Shalaby	mayshalaby@aucegypt.edu	Researcher A	3/2011-present
Amira Shoukry	am_magdy@aucegypt.edu	Researcher A	7/2011-present

Link Development Team

Amira Thabet	amira.thabet@mail.link.net	Researcher A, Team Leader	3/2011-present
Ashraf Hamed	ashraf.hamed@mail.lnik.net	Researcher and Developer	5/2011-present
Mohamed El Sherif	mohamed.elsherif@mail.link.net	Researcher and Developer	5/2011-present

The goal of this requirement specification document is to provide a unified view, common understanding of the functions of the prototype version for the Sentiment Analysis Tool for Arabic SATA and guidance to the development of the software system taking into consideration the project overall business plan, and user requirements.

Table of Contents

1. Introduction	3
1.1. Purpose	3
1.2. Intended Audience.....	3
1.3. Project Scope	3
2. General Description	4
2.1. Product Perspective	4
2.2. Product Functions	4
2.3. User Classes and Characteristics.....	5
2.4. Product Documentation.....	6
2.4.1. User Documentation.....	6
2.4.2. Technical Documentation	6
2.5. Product Constraints	7
3. Interface Requirements	7
3.1. User Interfaces	7
3.2. Hardware Interfaces	7
3.3. Communications Interfaces	8
3.4. Software Interfaces.....	8
4. System Features	11
5. Nonfunctional Requirements	13
5.1. Performance Requirements.....	13
5.2. Safety Requirements.....	13
5.3. Security and Privacy Requirements	13
5.4. Software Quality Attributes	14
6. Use Case Diagrams	14
7. Research Requirements	17

List of Figures

Figure 1: Architectural diagram for the product.....	5
Figure 3: Topic and Sentiment diagram	14
Figure 4: Hot Topics diagram	15
Figure 5: Sentiment diagram.....	15
Figure 6: Topic and Influential Users diagram	15
Figure 7: Sentiment and Influential Users diagram	16
Figure 8: Influential Users according to time.....	16
Figure 9: Influential Users diagram.....	17

1. Introduction

1.1. Purpose

This Software Requirements Specification (SRS) documents key specifications, describes a prototype in terms of functional and nonfunctional requirements for Sentiment Analysis Tool for Arabic (SATA). The information documented, helps the intended audience to design and develop the product. There will be a need for future updates of this document as we are planning to launch a prototype version for testing then start officially the beta version then the final version.

1.2. Intended Audience

Primary readers of this document are the web researches, web designers and developers. This document is intended for the following:

Developers: in order to be sure they are developing the right project that fulfills requirements provided in this document.

Testers: in order to have an exact list of the features and functions that has to respond according to requirements and provided diagrams.

Documentation writers: to know what features and in what way they have to explain. What technologies are required, how the system will response in each user's action etc.

System administrators: in order to know exactly what they have to expect from the system, right inputs and outputs and response in error situations.

1.3. Project Scope

The scope of the project is to provide a user friendly web based product that extracts people's sentiment feelings toward certain services, products, organizations, political or nonpolitical topics and any influential people on social media. In this project phase which aims at developing a filed prototype, emphasis will be put on Arabic tweets from Twitter in the political domain.

The project aims to:

1. Provide an accurate sentiment analysis results.
2. Achieve a wide range of users in Egypt and the MENA region.
3. Support Arabic Egyptian dialect in the first run and English will be considered later.
4. Smooth, fast, efficient, reliable and easy to use web-based tool.
5. Providing a user friendly menu and good entertainment visualization capabilities.
6. Having a plenty of options in term of filtering and viewing information according to user's needs.

2. General Description

2.1. Product Perspective

Due to the world's massive growth of social networks and the rapid flow of news over the internet; Link Development and AUC came up with the sentiment analysis tool for Arabic (SATA) research project. The main aim of SATA is that to develop a tool that can allow users to use a simple search bar to search for any services, products or any political topics and the engine of that tool is to crawl over the internet collecting all comments, reviews, tweets or even notes in blogs related to the user's search keyword. Then perform an intelligent processing technique to extract the true meanings of the people's comments and to decide and classify them in terms of positive, negative or neutral thus to know the majority of people like or dislike the desired topic. More specifically providing people's feelings regarding certain topics with high accuracy will lead to a better decision making.

The purpose of the prototype is to demonstrate the concept and to deliver operational and functional services for testing purposes. As for initial Twitter will be the only source of data for the prototype and then integration will be needed to include more sources like facebook, news websites and blogs.

2.2. Product Functions

The architecture diagram of the tool is shown in figure 1. This tool will provide the following functions:

- Topic Extraction

This part is considered a key stone in the project as it detects and extracts topics titles from the tweets. Using hash tags is not informative enough about the topic of sentiment the author mentions in his/her tweet. Our approach goes as follows, first we do preprocessing which includes removal of stop words that occur frequently in the tweets but have no relevant meaning, then generate the feature vector. The features used are n-grams, unigram, bigram, and trigrams, and some named entities that are extracted from the crawled tweets. The main step is to cluster related tweets together using similarity measures so we can have multiple clusters each has one topic. Afterwards key-phrase extraction is used on each cluster to extract the key-phrases that are candidates to be title topics. Clusters that result in similar key phrase are merged together and this key phrase has higher weight to be the topic title.

- Sentiment Classification

Sentiment classification is the primary module of the product. The objective of this part is to provide as much as possible an accurate classification for opinions embedded in certain sentences like tweets or micro-blogs written in Egyptian dialect as positive, negative or neutral. In addition to counting the total numbers of positive, negative and neutral tweets found in the data source with regards to specified topic.

- **Determining Influential Bloggers**

Since influential members in a social network can be responsible for starting a buzz or getting the community to notice a new trend, product, or even adopt an opinion, we are interested in the problem of identifying which users are leaders. For companies, organizations and governments, it is of great importance to learn about opinions in order to assess chances and risks. A manual analysis is only possible on a very limited scale. An automated computer supported analysis is necessary given the large number of virtual communities with huge amounts of postings.

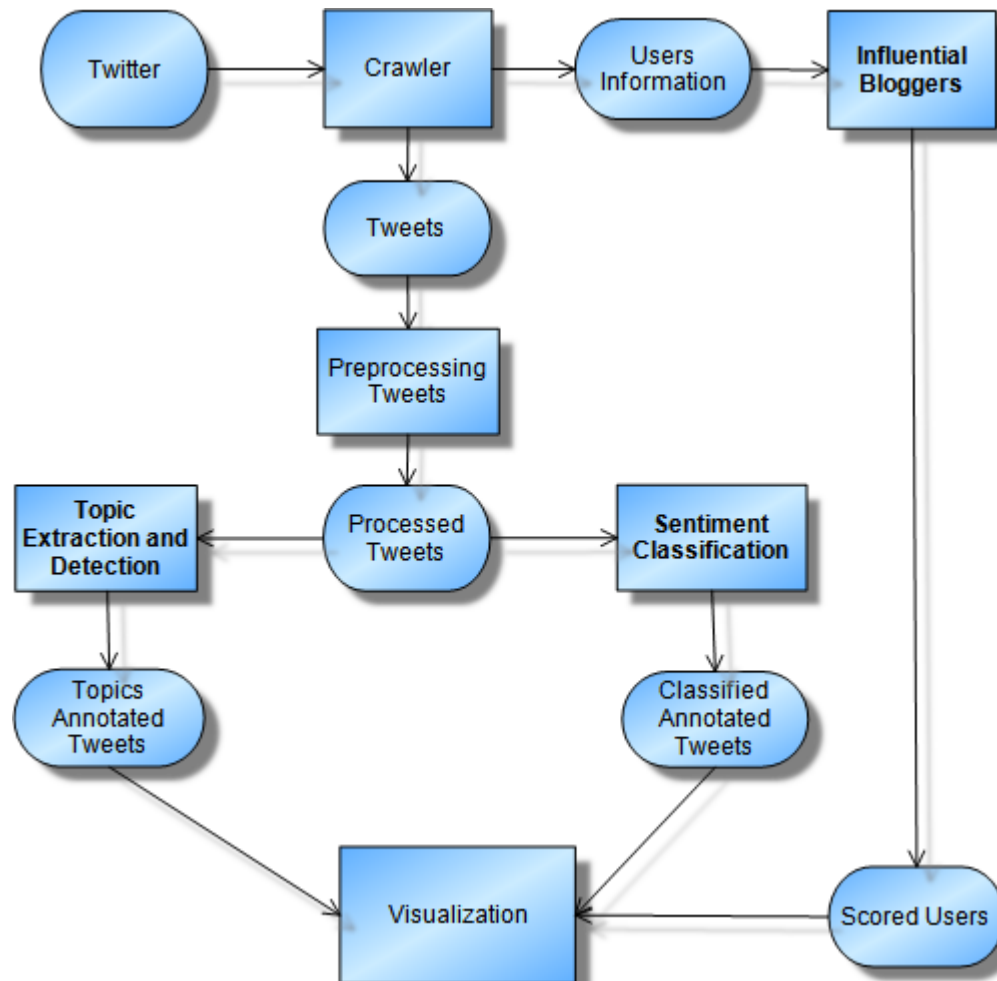


Figure 1: Architectural diagram for the product

2.3. User Classes and Characteristics

This part is to identify various user classes that we anticipate will use the web application. User classes will be differentiated based on the use, product functions and features, technical expertise, security and

privilege levels and educational level. The solution is intended to be used by three main different user classes; system administrators, system operators and customers or regular users.

No special knowledge or skills should be assumed for the part of the regular users. Users are not expected to learn or remember a set of commands in order to start using the application. The prototype application will be only a web based and then for the product versions there will be a desktop versions, smart phones and smart Tablets.

The following clearly describes a visionary role for each participant.

- **Users:** users with no particular knowledge needed, users who are interested to use the tool looking for knowing people's thoughts about a desired topic.
- **Advanced end users:** advanced users are those who have valuable input and feedbacks. Users who are more familiar with informative sites and can use our features efficiently. These valuable feeds will lead to enhancement of users' satisfaction.
- **System Operators:**
 - Maintains for the functional interface of the application and troubleshooting issues
 - Suggest possible updates and identifying renewal application needs
 - Coordinate with service providers and infrastructure vendors
 - Coordinate and communicate with system administrators
- **System Administrators:**
 - Develop and maintain installation and configuration procedures and operational requirements
 - Perform weekly/monthly backup operations, ensuring all required files and data are successfully backed up
 - Repair and recover from hardware or software failures
 - Coordinate and communicate with system operators

2.4. Product Documentation

2.4.1. User Documentation

User manual and CD will be made available for troubleshooting and help. Also this will represent as a full backup of the system. The user manual will contain detailed information about the usage of the product from a layman perspective to an advanced network/system administrator. The manual may also be made available online however this manual will be made for the product version but not for the prototype.

2.4.2. Technical Documentation

Technical manual will be made for the purpose of current and future developers involved in the product to understand and follow the solution at the level of coding and the programming languages used. The document will also include the development of technical requirements and the

functional specifications components for the sake of verifying the technical accuracy of all procedural steps included in the document to help in annual reviews process for developers over the product. Also as the user documentation this technical manual will be for the product version and not for the prototype.

2.5. Product Constraints

As we are planning to launch a prototype for testing purposes then a beta version for more advanced validation process then launching the final version. The following constraints will apply to for both the prototype and the different live service solution versions.

Processing Power:

SATA requires high speed machine for data capturing from various sources, classifying the sentiment polarity of large data and extracting topics.

Deployment Point:

SATA is built to be deployed as internet services. High bandwidth of the portal is required to fulfill the large number of concurrent users.

Operating Platform:

SATA may work for several distributions of Linux and Windows PCs, also smart phones and smart tablets.

3. Interface Requirements

3.1. User Interfaces

User interface includes various forms and windows. The main window will consist of the main search bar and a main menu bar with file, edit, view, tools and help. The interface will visualize the features and functionalities listed in this document for this prototype as the included below not limited to:

- Drop down menu for various option selection
- Selection list for filtering results
- Push buttons for users feedback and reclassifying tweets
- Visual graphs to show results
- Help button

3.2. Hardware Interfaces

The solution makes extensive use of several hardware devices. These devices include;

- MySQL database server with intensive use of memory space.
- PHP server with high performance and intensive use for CPU usage.
- Windows and Linux users' computers.

3.3. Communications Interfaces

Internet connection and a web browser are required in order to make use of several functions and to be executed such as searching, viewing and downloading.

3.4. Software Interfaces

For the prototype we will launch the portal over the internet and other than the hardware specified in the hardware interface section, the software requirements are to support windows operating system with support to MySQL, apache and PHP servers.

For the data gathering twitter is the only source and using Streaming API that offers high throughput. Using this API is perfect because we can retrieve real time information and also this continuous stream will be retrieved with no end and capturing all the messages in the stream without missing any information. The information retrieved in JSON format.

Twitter: Tweet

Basic information about a single tweet

No.	Name	Type	Contents
1	twitter.domains	Array of string	List of domains from links mentioned in this Tweet.
2	twitter.geo	Geo	The location from which this Tweet was sent.
3	twitter.in_reply_to_screen_name	String	The Twitter username of the user this Tweet is replying to if it is a reply.
4	twitter.links	Array of string	List of links mentioned in Tweet.
5	twitter.mentions	Array of string	List of Twitter usernames mentioned in this tweet.
6	twitter.source	String	The source of the Tweet. For example, "web" or "TweetMeme".
7	twitter.text	String	The text of the Tweet.

Twitter: User

Information about a user

No.	Name	Type	Contents
1	twitter.user.description	String	The Twitter user's biographical description.
2	twitter.user.followers_count	Integer	The number of followers the user has.
3	twitter.user.follower_ratio	Float	Ratio of followers to following users.
4	twitter.user.friends_count	Integer	The number of people the user follows.
5	twitter.user.id	Integer	Unique ID of the Twitter user.
6	twitter.user.lang	String	Two-character language code that the User set in Twitter.
7	twitter.user.listed_count	Integer	Number of lists in which the user appears.
8	twitter.user.location	String	The string description of the Twitter user's location.
9	twitter.user.name	String	The "real name" the user has assigned to themselves.
10	twitter.user.profile_age	Integer	The number of days since this user joined Twitter.

11	twitter.user.screen_name	String	The user's Twitter username.
12	twitter.user.statuses_count	Integer	The number of Tweets the Twitter user has posted.
13	twitter.user.time_zone	String	The Twitter user's time zone.
14	twitter.user.url	String	The URL the user added in their Twitter profile.

Twitter: User

Information about location

No.	Name	Type	Contents
1	twitter.place.attributes	Array of string	Additional information about the Twitterer's location.
2	twitter.place.country	String	The country from which this Tweet was sent.
3	twitter.place.country_code	String	Country code for the country this Tweet was sent from.
4	twitter.place.full_name	String	Full name of the place from which this Tweet was sent.
5	twitter.place.name	String	Short name of the place from which this Tweet was sent.
6	twitter.place.place_type	String	The type of place from which this Tweet was; for example: city, neighborhood, point of interest.
7	twitter.place.url	String	For a Tweet with place information, this string contains a link to the Twitter API to retrieve further information about the location.

Twitter: retweet

Information about tweet and the person who retweeted

No.	Name	Type	Contents
1	twitter.retweet.count	Integer	The total number of Retweets for this Tweet.
2	twitter.retweet.domains	Array of string	List of domains from links mentioned in the tweet that was Retweeted.
3	twitter.retweet.elapsed	Integer	In seconds how long between this retweet and the Tweet are they retweeting.
4	twitter.retweet.links	Array of string	List of links mentioned in the Tweet that was Retweeted.
5	twitter.retweet.source	String	The string source of the Retweet; for example: "web" or "Tweetdeck".
6	twitter.retweet.text	String	The Retweet text.
7	twitter.retweet.user.description	String	The biography information for the Twitter user who Retweeted this Tweet.
8	twitter.retweet.user.followers_count	Integer	The number of followers the user has.
9	twitter.retweet.user.follower_ratio	Float	Ratio of followers to following users.
10	twitter.retweet.user.friends_count	Integer	The number of people the Retweeting user follows.

11	twitter.retweet.user.id	String	The id of the Retweeting user.
12	twitter.retweet.user.lang	String	Two-character language code that the Retweeting user selected on Twitter's settings page.
13	twitter.retweet.user.listed_count	Integer	The number of lists the Retweeting user is listed in.
14	twitter.retweet.user.location	String	The string description of the Twitter's user's location.
15	twitter.retweet.user.name	String	The "real name" the Retweeting user supplied in Twitter's settings page.
16	twitter.retweet.user.profile_age	Integer	The number of days the Retweeting user been a member of Twitter.
17	twitter.retweet.user.screen_name	String	The Retweeting user's Twitter username.
18	twitter.retweet.user.statuses_count	Integer	The number of Tweets the Twitter user has posted.
19	twitter.retweet.user.time_zone	String	The Retweeting user's time zone.
20	twitter.retweet.user.url	String	The URL the retweeting user added in their Twitter profile.

Twitter: retweeted

Information about tweet and the person who retweeted

No.	Name	Type	Contents
1	twitter.retweeted.id	String	The unique ID of the Tweet that was Retweeted.
2	twitter.retweeted.mentions	Array of string	List of Twitter usernames mentioned in the Tweet that was Retweeted.
3	twitter.retweeted.place.country	String	The Retweeted Country from which this tweet was made.
4	twitter.retweeted.place.country_code	String	Country Code for the country this Retweeted Tweet was made from.
5	twitter.retweeted.place.full_name	String	Full name of the Place from which this Retweeted Tweet was made.
6	twitter.retweeted.place.name	String	Short name of the Place from which this Retweeted Tweet was made.
7	twitter.retweeted.place.place_type	String	The Retweeted type of place this tweet was made from.
8	twitter.retweeted.source	String	The source of the Retweeted Tweet. For example: "web" or "TweetDeck".
9	twitter.retweeted.user.description	String	The Retweeted Twitter user's description.
10	twitter.retweeted.user.followers_count	Integer	The number of followers the Retweeted author has.
11	twitter.retweeted.user.following_count	Float	Ratio of followers to following Retweeted users.

	ower_ratio		
12	twitter.retweeted.user.friends_count	Integer	The number of people the Retweeted user follows.
13	twitter.retweeted.user.id	String	The id of the Retweeted user
14	twitter.retweeted.user.lang	String	Two-character language code for the language the Retweeted user has set Twitter to.
15	twitter.retweeted.user.listed_count	Integer	Number of lists this Retweeted User is listed in.
16	twitter.retweeted.user.location	String	The string description of the Retweeted Twitter User has posted.
17	twitter.retweeted.user.name	String	The "real name" the Retweeted user has assigned to themselves
18	twitter.retweeted.user.profile_age	Integer	The number of days the Retweeted user has been a member of Twitter.
19	twitter.retweeted.user.screen_name	String	The Retweeted user's Twitter username.
20	twitter.retweeted.user.statuses_count	Integer	The number of Tweets the Retweeted Twitter user has posted.
21	twitter.retweeted.user.time_zone	String	The Retweeted user's time zone.
22	twitter.retweeted.user.url	String	The URL the retweeted user added in their Twitter profile.

4. System Features

This section illustrates the functional features using the following template:

System Feature: Name of the feature.

Priority: Indicate the priority of the feature to the user whether it is of High, Medium, or Low.

Description: Provide a short description of the feature

Action/ Response Sequences: List the sequences of actions required to be done in order to use this feature.

Result: List the system responses of this feature.

Functional Requirements: List the software modules required to carry out the function provided by the feature.

System Feature	Sentiment Classification
Priority	high
Description	Identifying the sentiment polarity (positive, negative or neutral) of tweets on certain topics from twitter.
Action	This module is activated after the user provides a query (topic, service or a product) or following the activation of the hot topic module.

Result	The system shows the results of the search of a query or the output of the hot topic module associated with the sentiment polarity of each item retrieved together with the percentage of Positive, Negative and Neutral sentiment of the whole result.
Functional requirements	A focused crawler, preprocessing module, sentiment classifier module, hot topic module and sentiment visualization module.

System Feature	User Feedback
Priority	Medium
Description	The user can give feedback by correcting the polarity of the classified retrieved tweets, and save the results
Action	The user selects a result and suggests a better annotation for it.
Result	The suggested correction by the user is stored in a system database to be handled by an administrator, and it is applied for future training and modifications to the system.
Functional requirements	A feedback interaction module

System Feature	Influential Bloggers Identification
Priority	medium
Description	Identifying the influential users on social media in certain topics.
Action	This module is activated after the user provides a query (topic, service or a product) or following the activation of the hot topic module.
Result	The system shows a list of all influential users on Twitter platform in certain topic, with indications on the level of influence.
Functional requirements	A focused crawler, Influential bloggers identification module, hot topic module, and influential blogger visualization module.

System Feature	Hot Topics Identification
Priority	High
Description	Identifying the Hot topics and Trending topics in Twitter according time period.
Action	This module is activated after the user provides a date interval. The default interval is the last week using the system date.
Result	The system shows the hot and trending topics, putting them in order from high trending topics to lower and the user can browse the tweets related to any of the topic.

Functional requirements	A focused Crawling, and the topic extraction module
-------------------------	---

System Feature	Results Visualization of the SATA components
Priority	medium
Description	Visualizing the results of sentiment classification, influential blogger and topic extraction modules into clear and interesting form.
Action	The proper modules will be activated by the user using a button included in the output screen of each of SATA modules.
Result	The system shows the results in the visualization form selected.
Functional requirements	Sentiment classification, influential blogger and topic extraction Visualization modules.

System Feature	Statistics and info-graphics
Priority	Low
Description	Viewing different collected statistics about retrieved classified tweets, hot topics tweets, and influential bloggers in a good visualized form such as info-graphics.
Action	The proper modules will be activated by the user using a button included in the output screen of each of SATA modules.
Result	Reports and Info-graphics that shows the statistics required
Functional requirements	Sentiment classification, influential blogger and topic extraction Statistics modules.

5. Nonfunctional Requirements

5.1. Performance Requirements

As for this prototype version we will keep on detecting if the system crashed, hanged or an operating system error occurred. Also detecting the performance of the system in terms of the efficiency of integration of the different components

5.2. Safety Requirements

For the safety requirements nothing but an operation of weekly backups for the data base should take place.

5.3. Security and Privacy Requirements

There are no specific security requirements, anyone can access and use the portal but only authorized persons who are allowed to use and access the database, web pages and the product engine.

5.4. Software Quality Attributes

- **Reliability**

The solution should provide reliability to the user that the product will run with all the features mentioned in this document are available and executing perfectly. It should be tested and debugged completely. All exceptions should be well handled.

- **Accuracy**

The solution should be able to reach the desired level of accuracy. But also keeping in mind that this prototype version is for proving the concept of the project.

6. Use Case Diagrams

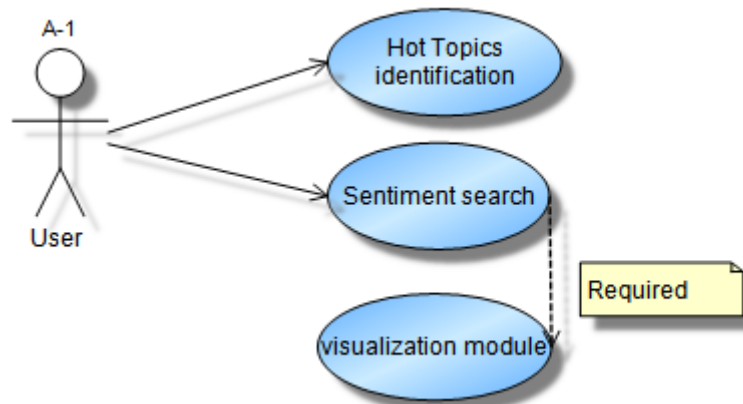


Figure 2: Topic and Sentiment diagram

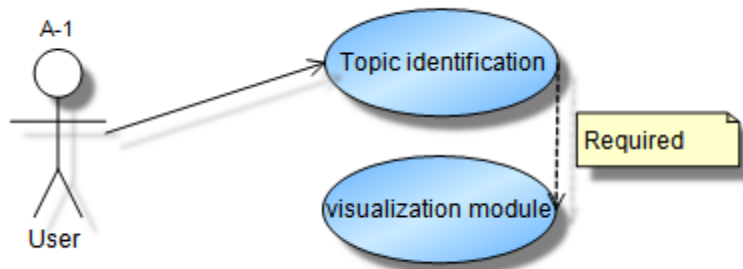


Figure 3: Hot Topics diagram

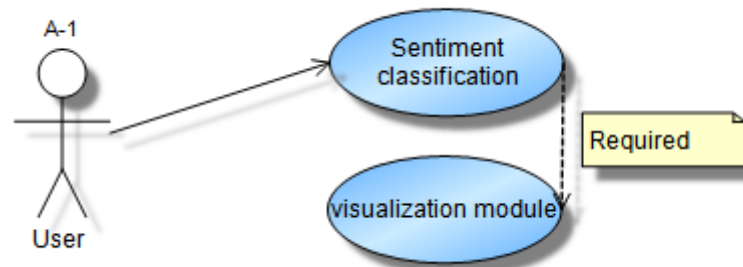


Figure 4: Sentiment diagram

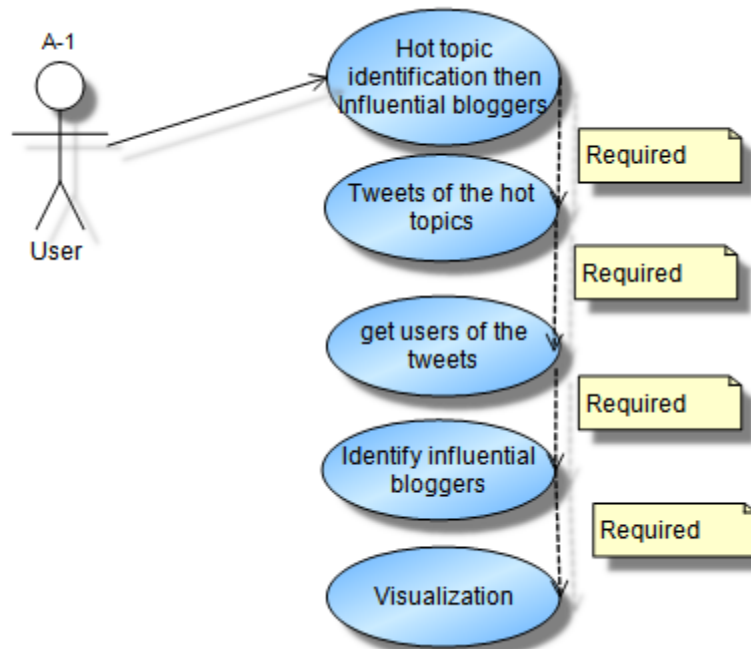


Figure 5: Topic and Influential Users diagram

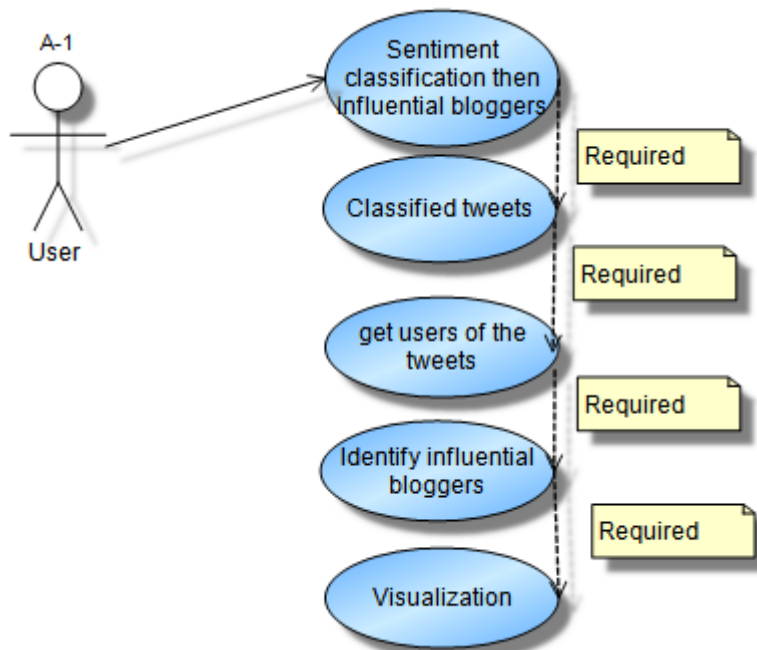


Figure 6: Sentiment and Influential Users diagram

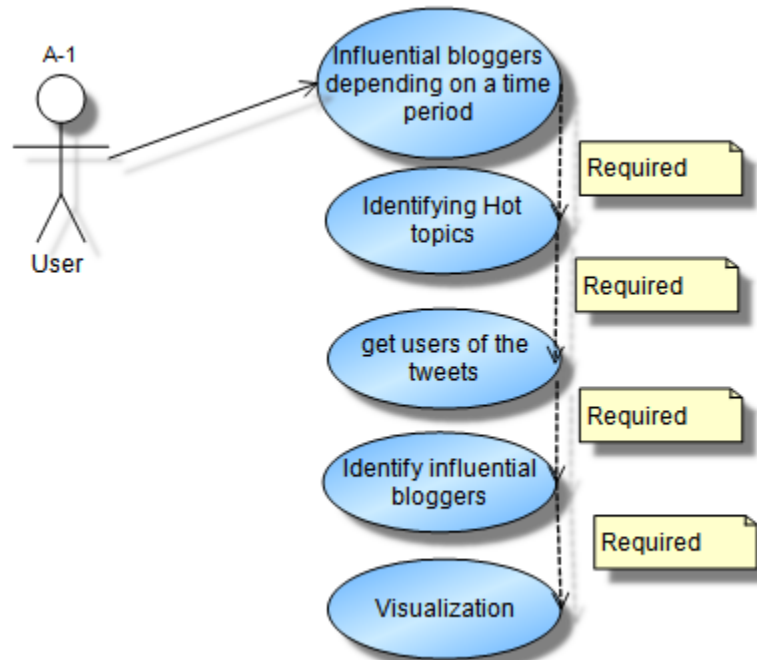


Figure 7: Influential Users according to time

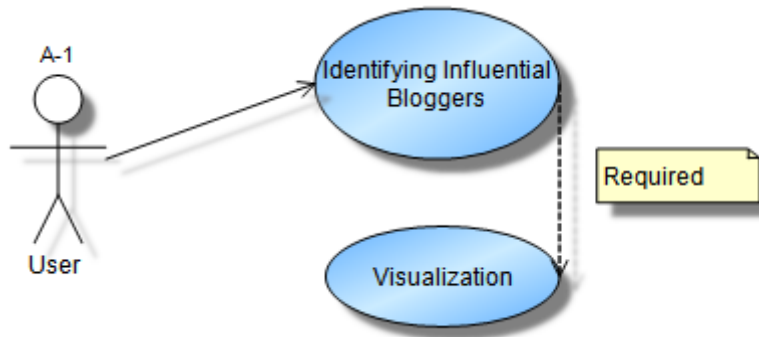


Figure 8: Influential Users diagram

7. Research Requirements

This section describes the needed research and experiments work efforts to develop each module: hot topic detection and extraction, sentiment classification, and detection of influential bloggers and opinion leaders.

Module(s) Name	Hot Topic Detection and Sentiment Classification
Research Objective	To find a list of Arabic stop words to be removed from tweets to enhance clustering and classification results.
Description	Finding a proper list of stop words is not an easy task specially when dealing with the Arabic dialect. Different spelling of the same word by users makes it difficult to include all the word forms in the list. Using natural language processing tools like stemmer to detect different forms of the same word is not just difficult but also gives bad results as some dialect words do not follow the inflection rules of modern standard Arabic. We will develop a list by getting frequent unigrams that occur more than a certain threshold from the total crawled tweets that reached about 20,000 tweets. Named entities are removed from this list as it's relevant to our work and being repeated that frequent gives it more weight not the opposite.
Expected Outcome	A list of Egyptian dialect stop words

Module Name	Hot Topic Detection
Research Objective	Select proper features that will achieve accurate clustering
Description	The features are the words or phrases that are relevant to our domain, and help clustering the tweets properly. In our work we are using n-grams and named entities. N-grams are unigram, bigram, and trigram. Determining the threshold of each n-gram that will lead to get better clustering is what we are targeting. The Named entities will be also considered as features will be also investigated.
Expected Outcome	Thresholds for all features

Module Name	Hot Topic Detection
Research Objective	Determine the most proper clustering technique that achieves high accuracy in grouping tweets that are talking about the same topic
Description	Different techniques are being used such as k-mean clustering; direct clustering, hierarchical, etc. The results of the clustering technique differ from one domain to another and from data set to another. It also differs according to the size of the data to be clustered. Some techniques works fine for small size of data and others works better for larger sizes. We will experiment with different techniques to decide on the best one using a large corpus of annotated data which contains pre-labeled tweets and use the clustering evaluation metrics : inter and intra similarity, entropy, and purity.
Expected Outcome	The best clustering technique to use

Module Name	Hot Topic Detection
Research Objective	Find the best way to label each cluster with the corresponding topic that suits it.
Description	<p>The approach that we will investigate is labeling by running key phrase extraction algorithm on each cluster to extract the main key phrases in it. A certain threshold is given to the algorithm so we can have a certain number of candidates. Choosing the best phrase from the resulted candidates is considered a problem. To solve this problem we will investigate these scenarios:</p> <ul style="list-style-type: none"> • If a named entity was discovered as one of the key phrases choose it as the most probable candidate • Use domain knowledge to choose the key phrase which is an entity in the ontology of the domain.
Expected Outcome	Accuracy of this labeling method

Module Name	Sentiment Classification
Research Objective	Extract the sentiment words in the tweets for the aim of creating a hybrid approach which combines the benefits of the ML approach and the SO approach.
Description	Given the limited work done for Arabic text in the field of sentiment analysis, especially for the Egyptian dialect, two lists of sentiment words will be built manually one for the most occurring positive sentiment words, and one for the most occurring negative sentiment words. Then for each word in these lists a weight is given to it based on its frequency in the positively labeled tweets, and the negatively labeled tweets in the corpus.
Expected Outcome	Weighted lists of the positive and the negative sentiment words mostly used by the Egyptian bloggers.

Module Name	Sentiment Classification
Research Objective	Determine the optimum threshold for each n-gram model separately, and for the combined n-grams models.
Description	A lot of studies had been made on the optimum threshold to use for each n-gram

	model (unigrams, bigrams, and trigrams) to present the text. These suggested thresholds are not necessary optimal in our case, which is presenting the tweets written in the Egyptian dialect. That is why each n-gram type will first needs to be tested separately using different thresholds in order to find the threshold which provides suitable cover to our sparse data. Second, these n-grams models will be combined together to further improve the performance of the classification process. Thus, for each combination, also the optimum threshold for each model will have to be figured out.
Expected Outcome	What are the suitable types of n-grams model to use in presenting the tweets, together with their optimum thresholds.

Module Name	Sentiment Classification
Research Objective	Compare the performance of the Machine Learning and the Semantic Orientation methodologies and choose the one which produces the best result.
Description	Although the ML approach was used extensively in the sentiment analysis process throughout the literature, it was still very important to test the SO approach with respect to our case which is dealing with the Egyptian dialect. Thus, we need to test both methodologies for the aim of comparing their performance and interpret the results obtained in each methodology.
Expected Outcome	The methodology which is most suitable to our case which is dealing with the Egyptian dialect.

Module Name	Choosing the Machine Learning classification algorithm
Research Objective	Compare the performance of the Support Vector Machine and the Naïve Bayes algorithms when used in the machine learning methodology and choose the algorithm which produces the best result.
Description	Although it was observed in more than one study that the Support Vector Machine algorithm produce higher result than the Naïve Bayes algorithm, it is still important to test the performance of the Naïve Bayes classification algorithm. The Support Vector Machine algorithm is believed to have some principle advantages over the Naïve Bayes algorithm. Some of these advantages are robustness in high dimensional spaces, any feature is relevant, robustness when there is a sparse set of samples and, finally, most text categorization problems are linearly separable. On the other hand, Naïve Bayes algorithms are also most suitable for classification problems with high dimensionality. That is why we need to try both algorithms and choose the one which produces the highest accuracy.
Expected Outcome	The classification algorithm which produces the highest accuracy.

Module Name	Detecting Influential Users and Opinion Leaders
Research Objective	Determine the method for retrieving User Information
Description	These two methods will be investigated: 1. Get the user information using the Twitter API

	<p>The Twitter REST API enables developers to access user information. However, the API is rate limited; it only allows clients to make a limited number of calls in a given hour. Also, there are limitations to the information retrieved, for example, it does not return more than 5000 followers per users even though the number of followers may exceed that, and for information such as retweets and mentions, it only returns the 20 most recent retweets or mentions, which may also exceed that. I find that such limitations may exclude information that could be of value to determine which users are influential.</p> <p>2. Get the user information from the user profile page Develop a crawler to access a user's twitter profile page source code, and retrieve the user information available, such as the number of tweets posted by that user, the number of followers, friends and list, and other information available that may be useful. However, there are a couple of issues in regard to this approach. First, Twitter has recently changed the layout of its user profile pages more than once. Such changes require adjustment to the crawler code that extracts the user information. Second, this approach limits the amount of information we may have access to what is only available on the page.</p>
Expected Outcome	Users information retrieval tool

Module Name	Detecting Influential Users and Opinion Leaders
Research Objective	Determine the Users' Scoring Model
Description	<p>Calculate the influence score of a user, excluding news site from the scoring since they wouldn't be taken in consideration as influential members. From a sample of 19880 tweets retrieved using the query "tahrir", 650 users were extracted after excluding news site. The data collected was very sparse as can be seen in the following values of the parameters that can be used to score the influence of a user:</p> <ol style="list-style-type: none"> 1. Total Number of tweets: 5..109508 (<i>zeinobia</i>) 2. Number of followers: 0..507946(<i>el baradei</i>) 3. Number of Friends: 2..3845 4. Number of Listed 0..7624 (<i>el baradei</i>) 5. Number of mentions and retweets in the sample collected 0..962 6. Number of tweets for a user in the sample collected 1..790 <p>Given the nature of the data, we want to find a way to use it to our advantage and develop a model for scoring the user influence.</p>
Expected Outcome	An equation that uses the user information and data to produce a score that to rank the users based on how influential that user is.

Module Name	Detecting Influential Users and Opinion Leaders
Research Objective	Decide on the Evaluation Method
Description	<p>There is not training or testing data to evaluate the efficiency of any proposed model, nor is there an obvious reference point with accurate information regarding influential users on twitter. For the literature review, one paper resorted to the website Digg (http://www.digg.com) to provide a reference point. Digg is about user powered content, and can be considered a large online user survey, since everything is voted on, however that may not be applied to Twitter members' posts. Another paper which was studying influence on Twitter, studied only URL tweets so that I could use Bit.ly (https://bitly.com/) which is a URL shortening service that for each shortened URL keeps track of how many times it has been accessed, so the bit.ly URLs found in tweets can be queried for the number of clicks the service has registered on that URL. The URL click data was used to test how well the influence measure can predict the attention the URLs posted by the users receive. That, however, may not be applied in our case either. We will probably resort to a manual evaluation, but with the assistance of any available online reference points about the user in question. From the user's Twitter profile page we may get a glimpse of who that person is or for example, we may be able to find the user's facebook page, or we could search Google to find any online record or information, web sites, citations, CV...etc. In short, any additional information to assist the evaluation and help us determine whether that user is in fact a person with influence or not.</p>
Expected Outcome	Manual evaluation method, and reference point, to confirm that our scoring model is producing results that are accurate; that the high ranking users are in fact people who may be considered leaders or influential members in a community, in regard to a certain topic.