# Feature-Based Sentiment Analysis for Arabic Language

Eng. Ghady Alhamad[1], Dr. Mohamad-Bassam Kurdy[2]

Master in Web Science, Syrian Virtual University, Hama, Syria[1]

Ph.D. in Mathematical Morphology, Syrian Virtual University, Dijon, France[2]

*Abstract*—**In light of the spread of e-commerce and e-marketing, and the presence of a huge number of reviews and texts written by people to share views on products, it became necessary to give attention to extracting these opinions automatically and analyzing the feelings of the reviewers. The goal is to obtain reports evaluating products and contribute to improve services at a glance. Sentiment Analysis is a relatively recent study that deals with the processing of natural texts published in web sites and social networks. However, the processing of texts written in the Arabic language is one of the challenges that specialists face because people do not rely on standard Arabic, writing people in spoken/colloquial languages and use various dialects. This paper will present feature-based sentiment analysis for Arabic language which works on text analysis technique that breaks down text into aspects (attributes or components of a product or service), and then allocates each one a sentiment level (positive, negative or neutral).**

*Keywords*—*Sentiment analysis; feature-based; colloquial Arabic; opinion mining; natural language processing*

## I. INTRODUCTION

Sentiment analysis is an active research area since 2003 [1] and, it refers to the process of mining the texts in order to identify the tone of the passage written by the reviewers [2]. These tones are the focus for the decision makers to assess customer satisfaction with their products, which have been categorized into different poles. The most significant polarizations were absolutely in many studies, such as [3], [4] and others were usually three tones: positive, negative, and neutral. Sentiment analysis, which is also called opinion mining, is the computational study of people's opinions, sentiments, and attitudes about topics, entities, people and events, that are expressed in texts [5].

Recently the number of internet users has increased significantly in the Middle East and people are becoming more and more interested in buying online. According to new statistics [7] which have resulted that the number of internet users in the Arab countries has reached 157 million people, according to the Arabic Network for Human Rights Information. Internet buyers are distributed in the Middle East in several countries, reaching 10.6 million in Saudi Arabia, 6.8 million in the UAE, 2.4 million in Kuwait and 15.2 million in Egypt and around other Arab countries at different rates. The mobile phone is also the best-selling product online in the Arab world, according to the director of Souq. (Source: payfort) [7].

The Arabic language is one of the fastest growing languages on the web [6]. The main challenge in this study that sentiment analysis is for Arabic which is considered a poor area for this language. In addition to the peculiarity of the Arabic language whether in the Standard Arabic or in terms of the diversity of its dialects. The Arabic language is a Sematic language which consists of 28 letters. It is a cursive language, in which word formation consists of connecting letters to each other. As opposed to the English language. Arabic writing starts from right to left and has no capitalization [6].

Human can easily read texts and recognize reviewer's sentiment by understanding context, but for computers it is not normal process. Therefore, the main task in this study is to make computers recognizing the reviewer's sentiment and this achieved by Natural Language Processing (NLP). NLP is a framework to support an interaction between computers and human languages [8].

In this paper, based on the market need in the Arab world, and in light of the lack of Arab studies in this field with the wide spread of Arabic texts on the web written in various non-standard Arabic dialects, it was necessary to fill the gap and present a theory in this field. Since mobiles are the best-selling products, they will be the focus of this study. This theory exhibits a proposed method for recognizing Arabic sentiment phrases for mobile phones with consideration of each feature of phone like: camera, battery, memory … etc. The opinion phrases identified by building grammatical analyzer which is defining several forms for these phrases. Grammatical analyzer needs a lexical analyzer as input to define opinion tokens. opinion tokens could be mobile features, entities names and opinion words. Which could be positive names, negatives names, positive verbs, negative verbs, positive adjectives, negative adjectives, modifiers and negation words. This process called Parts of Speech Tagging (POS) that will be presented in this study. POS tagging has been used for a long time in text classification and NLP. POS tagging differentiates syntactic meaning of words in a sentence by using some specific tags, such as tags for noun, pronoun, verb, adjective, adverb, conjunction and others [8].

Also, after identifying opinion phrases the study will classify the opinion into five polarities in range [-2, 2]: {Strong positive, positive, neutral, negative, strong negative}. Finally, the summarization is necessary in order for decision-makers to gain knowledge.

The rest of this paper is organized as follows. Section 2 overviews related work. Section 3 describes the methodology followed with examples showing exactly how the study could achieve the goal. Section 4 presents the results of the experimental analysis and evaluation. Finally, in section 5, conclusions and possible future work are discussed.

## II. RELATED WORKS

This section exhibits a number of related previous studies as this paper adopts some of their approaches and overcome the absence of some points for Arabic in others.

Bing Liu [9] is one of the most famous studies that cited by most researches in this field. He used rules for recognizing opinions and made for them [1]Backus-Naur Form (BNF). BNF is a meta syntax notation for context-free grammars, often used to describe the syntax of languages used in computing. This study adopted his approach for Arabic language. In Mohammad N. et al. [1] recognized opinions using lexicon, they concerned in Modern Standard Arabic (MSA) and colloquial for example: "Khaliji". In addition to Chetashri B. et al. [10] discussed the lexical and machine learning approach. Mongkol Seansuk et. al. [11] exhibited traditional methodology and evaluated opinions logically for each sentence; they considered opinion is positive by comparing sentences and the result will be positive only if both of them are positive, else negative.

Asad Ullah R. K. et al. [12] retrieved comments from YouTube to analysis sentiment about Android and iOS; they used General Architecture for Text Engineering ([2]GATE) component and build plugin. GATE is necessary component for NLP, as this paper used it to achieve multiple ideas. Weishu Hu. et al. [13] presented how to mine product features in opinion sentences. It made use of SentiWordNet based algorithm to find opinion of the sentence. Samir A et al. [14] presented a novel solution for Arabic Named Entity Recognition (ANER) problem, which aimed to boost the identification of extracted named entities. They utilized a machine learning technique using pattern recognition to classify name entities (NE).

Sana A. et al. [6] proposed study for Twitter sentiment analysis model that based on supervised machine learning and semantic analysis. They are divided their approach to two phases training and testing, in the training phase, they needed to learn from a set of labeled tweets for classifier. Then they used to classify unlabeled tweets in the testing phase. Mohamad H. et al. [16] also focused on studying sentiment analysis for Arabic language that collected from Twitter, Facebook and YouTube. Taysir .H et al. [15] focused on mining social networks for sentiment analysis of colloquial Arabic comments. The approach concerned with Egyptian

terminology as it provided a structure to define the standard meaning of the word and the informal terms associated with this word. Alaa El-Dine A. H. et al. [17] also used classification methods to analyze users' comments and detected the comments that agree, disagree or is natural with respect to a post. The data collected from Facebook.

Sawsan C. et al. [18] adopted in their approach ontology for detecting Arabic Emotion. They detected language or dialect that belonged to with the help of GATE. They arrange the emotional vocabularies into intensities belonging to the integer numerical domain [-10, +10]. Whereas other studies detected specific dialect of Arabic language like Abdullah D. et al. [19]. Arabic Levantine tweets are a corpus of the study, they implemented different methods to automatically classify text messages of individuals to infer their emotional states.

Abdul-Mageed et al. [20] presented a subjectivity and sentiment analysis system (SAMAR) based on a Support Vector Machine (SVM) classifier for different Arabic social media applications: Web forums, chat, Wikipedia Talk Pages, and Twitter. They studied different features including word n-grams, POS tagging, and word stems. Also, many stylistic features related to social media applications were investigated. The results showed that the classifier performance relied on the type of the dataset and feature used.

## III. METHODOLOGY

This section presents method for feature-based sentiment analysis for Arabic language. Mobile phone is the target product. Therefore, the study exhibits analyzing people's sentiment for mobile phones for each feature. As well as it presents entity recognition for mobile names. The study consists of the process as shown in Fig. 1.
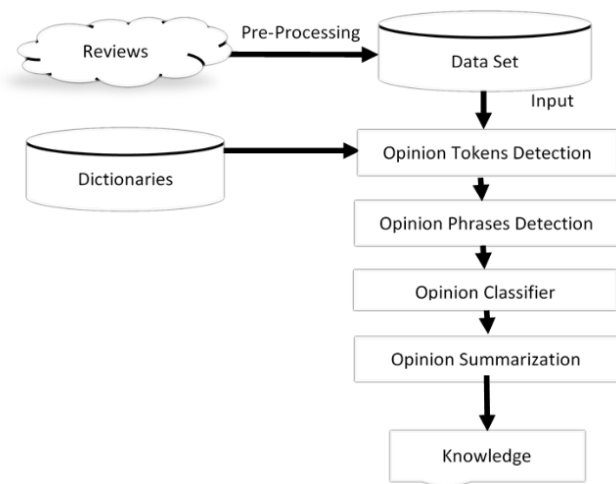


Fig. 1. Feature-based Sentiment Analysis for Arabic Language.

[1] *https://en.wikipedia.org/wiki/Backus%E2%80%93Naur_form*

[2] *https://gate.ac.uk/*

## A. Dictionaries

This approach includes three dictionaries for features, sentiment words and entities. The data collected as a sample for training data. With the possibility of feeding these dictionaries later dynamically with flexibility, prior experience is not required. The dictionaries data collection source details in pre-processing section.

*1) Features dictionary:* Features are domain-based of sentiment analyzer and in this study the domain is mobile phones. This dictionary composes 84 words as a sample data and it is scalable.

E.g. for mobile features: camera; "كاميرا", memory; "ذاكرة", battery; "البطارية" ...etc.

*2) Sentiment words dictionary:* Sentiment words contribute to the quality of sentiment classifier. They are domain-independent unlike the features, but they are related to the terminology of the Arabic language in all its dialects. They were collected by relying on experiments from people's reviews, and space was also allowed for scalable.

Sentiment words are classified into several categories: five positive categories, five negative categories, negations category, and strong words (or modifiers) category.

The structure of sentiment words is shown in Table I. with number of words for each category and examples.

TABLE I.  SENTIMENT WORDS STRUCTURE

| Sentiment Words Categories | | |
|---|---|---|
| *Positive* | *words* | *Example* |
| Strong Positive Adjectives | 25 | wonderful - "رائع" |
| Positive Adjectives | 46 | nice - "حلو" |
| Positive Names | 22 | advantage - "ميزة" |
| Positive Verbs | 19 | advice - "أنصح" |
| Positive Comparitive | 13 | better - "أفضل", "أحسن" |
| *Negative* | *words* | *example* |
| Strong Negative Adjectives | 7 | disappointed - "مخيب" |
| Negative Adjectives | 23 | weak - "ضعيف" |
| Negative Names | 23 | disadvantage - "عيب" problem - "مشكلة" |
| Negative Verbs | 20 | hate - "أكره" |
| Negative Comparitive | 8 | worst - "أسوأ" |
| *Other* | *words* | *example* |
| Negations | 12 | not - "ليس", "ما" |
| Strong Words (modifiers) | 13 | very - "بزاف", "وايد", "كثير" |

*3) Entities dictionary:* Entities represent product names, which was one of the challenges this approach really face. Because there is no standard way to write mobiles names. Most of the mobile phone brands are not Arab. The main problem is when someone tries to write mobile name in Arabic alphabet. In addition, it may not be strange if others mixed spelling the name between Arabic and Latin letters.

An example of the different cases that reviewers write for the mobile name of the "Samsung Galaxy S6".

*a)* جسامسونجS6. (Mixed, not full name).

*b)* Samsung S6. (Review in Arabic, mobile name in English, not full name).

*c)* جلاكسي اس6. (Arabic only, not full name with other parts).

*d)* جلاكسي S6. (Mixed, Not full name).

*e)* S6. (Version only, not Arabic).

*f)* اس6. (Version only, Arabic).

Also, the same word may have several spellings in Arabic, that the stemmer unable to stem them because that are not Arabic and have no meaning.

- "غلاكسي", "جلكسي", "جالاكسي", "جلاكسي", (for Galaxy word).

- "سامسونغ", "سامسونج", (for Samsung word).

Therefore, this approach defines specific structure for mobile names as a hierarchy. Each level has multiple keywords to include all different spellings for the same name.

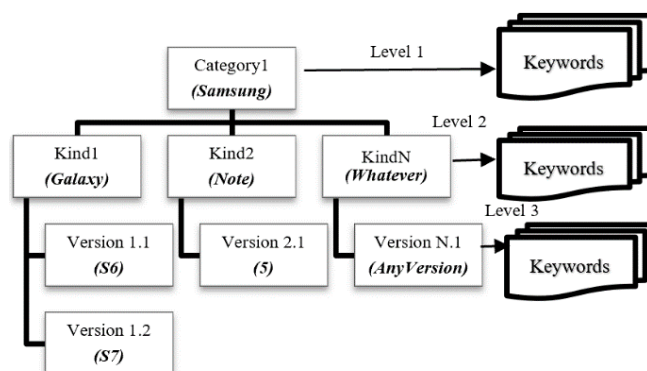Fig. 2 shows the entity structure categorized into three levels with examples.



Fig. 2.  Entity Data Structure.

This dictionary includes three mobile categories (brands), each category has several kinds and each kind has several versions. Each level has keywords list for different forms of the same word.

## B. Pre-processing

*1) Data collection:* As sentiment analysis depends on the training data which labelled. The 85 posts of mobile phones are collected as a dataset, they include 1024 comments, which include 570 replies obtaining from mobile pages; [3]souq.com and [4] mobihall.com pages on Facebook. Most of posts are advertisement about mobiles therefore, the comments and replies are the target reviews.

*2) Reviews structure and format:* Since the reviews has been collected from different sources, the standard structure

---

[3]*http://souq.com*
[4]*http://mobihall.com*

became required. The appropriate format chosen to represent reviews is the "eXtensible Markup Language" (XML). The reviews include ratings as likes or stars, created date time and review id. Fig. 3 shows sample of reviews data.

```
<posts>
<post>
<id>261</id>
<url>mobihallsocial</url>
<socialmedia>Facebook</socialmedia>
<postisadvertisment>false</postisadvertisment>
<createdtime>16/04/2016</createdtime>
<eval>
<value>25.0</value>
<type>LIKES</type>
</eval>
<message> سامسونج جلاكسى اس 6 موبايل انيق والمعالج الاكسيون اثبت
كفائة عالية ال3 جيجا بايت ليست كافية ابدا الكاميرا ممتازة تصميم اكثر من رائع
العيوب اللى قابلتها عدم وجود فتحة لكارت الميمورى وشريحة واحد
اتصال</message>
<comments>
<comment>
<id>2763</id>
<createdtime>16/04/2016</createdtime>
<eval>
<value>10.0</value>
<type>LIKES</type>
</eval>
<message>وفي عيب اخر وهو البطاريه</message>
<replies>
<reply>
         <id>585</id>
         <createdtime>16/04/2016</createdtime>
         <eval>
         <value>0.0</value>
         <type>LIKES</type>
         </eval>
         <message> انا عندى نفس المشكله????</message>
</reply>
</replies>
</comment>
</comments>
</post>
</posts>
```

Fig. 3. Sample of Reviews Sata.

As well as the dictionaries data that collected by reliance on the same sites, they are not only for reviews that represent dataset but also for all reviews of whole sites pages.

*3) Arabic stemmers:* This approach needs two stemmers for reviews and dictionaries lists:

*a) Light stemmer* is built-in by this study for noise elimination or normalization:

✓ Standardize Hamza "أ".

✓ Eliminate Tashkeel ٓ, ٔ, ٰ, etc.

✓ Standardize "ة", "ه".

✓ Standardize "ي", "ى" in the end of the word.

✓ Remove repeated letters "جداااااااا".

✓ etc.

*b) Advance stemmer* to extract root words. "Khoja" and "Arnlp" are the most famous stemmers for Arabic language.

This approch used "Arnlp" because in addition to finding the root word, it works to find the stem word. The stem word may be more meaningful and reduce the confusion that occurs due to the presnce of one root for opposing words.

*C. Sentiment Analyzer*

The approach achieves natural language processing with GATE component. The process of sentiment analyzing consists of these steps:

*1) Opinion tokens detection:* The detection of opinion tokens is considered as the lexicon in this study. Opinion tokens include dictionaries lists; features, sentiment words and entities. The opinion tokens detection implemented using GATE Gazetteer. The GATE Gazetteer matches words in lists with the possibility of annotating each matched word. These annotations are very useful data for next steps. Therefore, Gazetteer includes following lists:

*a) Features list includes:* feature, feature identifier.

*b) Sentiment words list includes:* sentiment word, sentiment category, polarity, sentiment id..

*c) Entities list for all keywords of mobiles names in one list, includes:* keyword of mobile name, full mobile name, product id, level name.

Fig. 4 shows opinion tokens detection example. The lexical analyzer detects opinion tokens and classifies them based on its semantic meaning as kind of POS tagging.



Fig. 4. Opinion Tokens Detection Example.

*2) Opinion phrases detection:* In this stage, the grammar analyzer is identifying opinion phrases syntax. The opinion phrases detection is performed using GATE JAPE transducer. The GATE JAPE transducer defines rules for forms of opinion phrases. It takes opinion tokens annotation as input and the output is the opinion phrases annotation. Before identifying the opinion phrases rules, using JAPE the reviews should be split into sentences by GATE Sentence Splitter to detect phrases for each sentence separately. Opinion phrases rules described with [1]BNF meta syntax notation that used to describe the syntax of phrases.

Fig. 5 shows the suggested syntax by this approach of opinion phrases for Arabic language. The BNF identifies six rules (or cases) for the forms of opinion phrases rules, and two rules for compare two products.

```
<OpinionRule_1> ::= (Product | Feature)? (Negation)?
(SentimentWord) (Modifier)?

<OpinionRule_2> ::= (Feature) (Product) (Negation)?
(SentimentWord) (Modifier)?

<OpinionRule_3> ::= (Feature | Product)(((Modifier)?
(SentimentWord)) | ((Negation)? (SentimentWord)
(Modifier)?))+

<OpinionRule_4> ::= (Feature) (Product) (((Modifier)?
(SentimentWord)) | ((Negation)? (SentimentWord)
(Modifier)?))+

<OpinionRule_5> := (Negation)?
     (((Modifier) (SentimentWord)) | ((SentimentWord)
     (Modifier)) | (SentimentWord)))
  ((Product ) | (Feature) (Modifier)?)+

<OpinionRule_6> ::= (Negation)? (((Modifier)
(SentimentWord)) | ((SentimentWord) (Modifier)) |
(SentimentWord)) ((Feature) (Product))+

<CompareOpinionRule_1> ::= (Feature1)? ( Product1 |
Category1 | Kind1) (Negation)? (SentimentWord: Comparative)
(Feature2)? ( Product2 | Category2 | Kind2)

<CompareOpinionRule_2> ::= ( Product1 | Category1 | Kind1)
(Negation)? (SentimentWord: Comparative) (Feature1)
```

Fig. 5. Opinions Phrases BNF.

Notes about BNF:

*a)* Product is the entity with full mobile name.

*b)* Category is the first part of mobile name, that often represents company name.

*c)* Kind is the second part of mobile name.

*d)* Version is the third part of mobile name.

Fig. 6 shows opinion phrases detection example. The grammatical analyzer detects opinion phrases to test eight rules. The bottom table in the figure shows necessary information, opinion words, polarity, mobile features, mobile name and rule identifier.



Fig. 6. Opinion Phrases Detection Example.

*3) Opinion classifier:* This stage classifies opinions into polarities. The polarities divided into five categories; {Strong Positive, Positive, Neutral, Negative, Strong Negative}.

The classification begins for each of the opinion phrases that are defined in the BNF. The opinion phrase cannot be neutral at all, but it is possible that the entire review is neutral when there is balance between the positive and negative opinion phrases in the same review. The polarities have been defined mathematically in the following ranges values between [-2 and +2]:

*a) Stron*g Negative: [-2, -1.5].

*b)* Negative: ] -1.5, -0.5].

*c)* Neutral: ]-0.5, 0.5].

*d)* Positive: ] 0.5, 1.5].

*e)* Strong Positive: ] 1.5, 2.0]

Since the review consists of one or more opinion phrases and polarity value is fuzzy in the range [-2, +2], the polarity will be calculated by average function for polarities of opinion phrases.

$$\textbf{\textit{Polarity}} = \frac{\sum Polarity_{opinion\_phrase} \times \left(weight+1\right)}{\left(weight+1\right)} \tag{1}$$

Weight either represents the number of opinions for specific polarity or the number of likes. Likes mean if someone has copied the same opinion and gets the same polarity. The polarity value is multiplied by weight + 1. +1 represents the opinion itself and avoids dividing by zero.

Suppose 9 positive opinion phrases, 1 strong positive opinion phrase, 1 negative opinion phrase, and 4 strong negative opinion phrases. Where Strong Positive 2, Positive 1, Negative -1, Strong Negative -2. The (2) shows for polarity calculation.

$$\textbf{\textit{Polarity}} = \frac{9\times\left(1\right)+1\times\left(2\right)+1\times\left(-1\right)+4\times\left(-2\right)}{9+1+1+4} = \frac{2}{15} = 0.1333 \tag{2}$$

Final result based on ranges that are shown in the beginning of this section where $0.1333 \in ]-0.5, 0.5]$, therefore, it is Neutral.

The Table II shows example for each rule defined in the BNF. It must be pointed out that in Arabic grammar, the noun comes before the adjective, in contrast to English grammar, where adjective precedes the noun that is being described, in addition to some other differences, therefore the translation of examples is only for illustration and it is not necessary that is correct for English. For example, "Red Flower", in Arabic, it is written as "Flower Red" – "الزهرة حمراء". Therefore, the illustration respects BNF rules and word order.

TABLE II.  EXAMPLE FOR BNF RULES (OPINION PHRASES)

| Opinion Phrases Examples | | |
|---|---|---|
| *Rule Name* | *Example* | *Polarity* |
| OpinionRule_1 | سامسونج جلاكسي اس6 جميل جداً / Samsung Galaxy S6 is nice so much | 2 |
| OpinionRule_2 | كاميرا السامسونج جيدة جداً / Camera of Samsung is good very much | 2 |
| OpinionRule_3 | جـ يد 6 اس أي فون ومم يز / Iphone 6s is good and distinctive. | {1, 2} |
| OpinionRule_4 | كاميرا أيفون اس6 رائعة جداً ومميزة. / Camera of iPhone 6s is wonderful very much and distinctive | {2, 2} |
| OpinionRule_5 | لا أنصح بالسامسونج اس7 و اس6 / I do not advise you with Samsung s7 and s6 | -1 |
| OpinionRule_6 | لا أنصح بكاميرا سامسونج اس6 / I do not advise you with camera of Samsung s6. | -1 |
| CompareOpinionRule_1 | كاميرا سامسونج اس 7 أفضل من كاميرا سامسونج اس 6 / Camera of Samsung s7 is better than camera of Samsung s6 | 1 for product1 -1 for product2 |
| CompareOpinionRule_2 | أفضل جهاز بالعالم / The best device in the world | 1 |

*4) Opinions summarization:* This method summarizes results of sentiment analyzer in several ways:

*a) Feature-based:* The method summarizes the results of opinion polarity for each feature separately.

*b) ReviewDate-based:* The method summarizes the results of opinion polarity during specific time periods.

*c) Polarity-based:* The method summarizes the results of opinion polarity as a ratio for each polarity.

*d) Product-based:* The method summarizes the results of all opinions polarity.

## IV. EXPERIMENTS AND RESULTS

This section presents the experiments and results. The experiments are performed to analyze the quality of the proposed methodology whereas in the results will present the results of this study with examples. The experiments are achieved with precision, recall and f-measure for opinion tokens detection and opinion phrases detection.

As for opinion tokens detection are evaluated by the dictionaries size and by stemmer for matching words. The Table III shows a test with 20 tokens extracted from several reviews consists of 100 words.

TABLE III.  OPINION TOKENS DETECTION EXPERIMENT

| | *Positive (Retrieved)* | *Negative (Not Retrieved)* |
|---|---|---|
| *True* | 18 | 77 |
| *False* | 3 | 2 |

The Precision for opinion tokens is defined by the formula:

$$P = \frac{TP}{TP+FP} = \frac{18}{18+3} \approx 0.86 \tag{3}$$

The Recall for opinion tokens is defined by the formula:

$$R = \frac{TP}{TP+FN} = \frac{18}{18+2} \approx 0.9 \tag{4}$$

From (3) and (4) results, the F-measure is defined by the formula:

$$F = 2.\frac{R \times P}{R+P} = \frac{2 \times (0.9 \times 0.86)}{0.9+0.86} \approx 0.88 \tag{5}$$

As well as the opinion phrases detection are evaluated by measuring the quality of opinion rules that defined in BNF in the METHODOLOGY section. The Table IV shows a test with 35 phrases extracted from several reviews consists of 60 phrases.

TABLE IV.  OPINION PHRASES DETECTION EXPERIMENT

| | *Positive (Retrieved)* | *Negative (Not Retrieved)* |
|---|---|---|
| *True* | 31 | 20 |
| *False* | 5 | 4 |

The Precision for opinion phrases is defined by the formula:

$$P = \frac{TP}{TP+FP} = \frac{31}{31+5} \approx 0.86 \tag{6}$$

The Recall for opinion phrases is defined by the formula:

$$R = \frac{TP}{TP+FN} = \frac{31}{31+4} \approx 0.89 \tag{7}$$

From (6) and (7) results, the F-measure is defined by the formula:

$$F = 2.\frac{R \times P}{R+P} = \frac{2 \times (0.89 \times 0.86)}{0.89+0.86} \approx 0.87 \tag{8}$$

The results show in followed figures some models for opinions summarization to build knowledge that can benefit decision makers: Fig. 7 shows bar chart for feature-based statistics about comparison of two mobiles Sony Xperia Z5 and Sony Xperia Z3. It shows polarity for each feature in range [-2, 2].
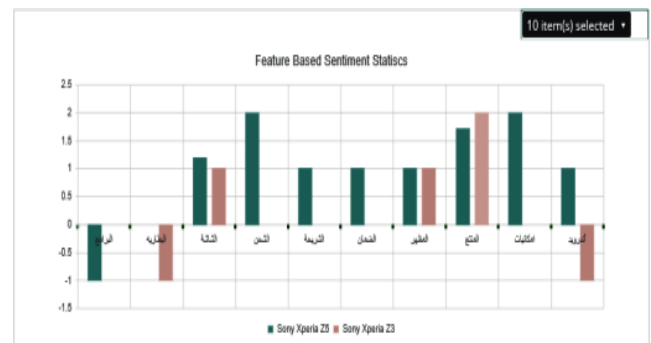


Fig. 7.  Feature-Based Summarization.

Fig. 8 shows pie chart that summarizes the percentage rate for each polarity.
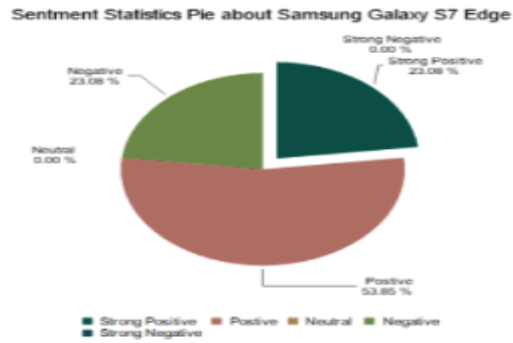


Fig. 8. Polarity-Based Summarization.

Fig. 9 shows line chart about comparison of two mobiles that summarizes the polarity for each mobile through specific periods. It illustrates statistics from 2014 to 2016. Polarity in range [-2,2] for mobiles Apple iPhone 5s and Sony Xperia Z5.



Fig. 9. ReviewDate-Based Summarization.

Fig. 10 shows gauge chart about comparison of two mobiles that summarizes the final polarity for each mobile. Polarity in range [-2,2] for mobiles Apple iPhone 5s and Sony Xperia Z5.
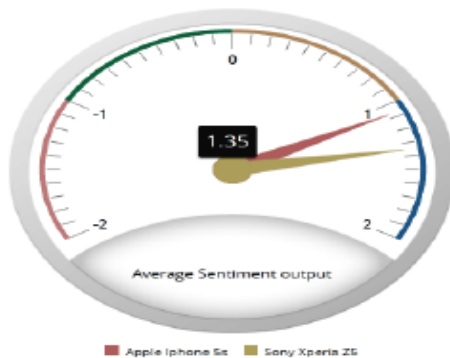


Fig. 10. Product-Based Summarization.

## V. Conclusion

This paper proposes feature-based sentiment analysis for Arabic language, the target product is mobile phone. The results of this mining are demonstrated as the degree of strong positive, positive, neutral, negative and strong negative. This result is useful for both consumers and companies. This study presents an approach in active area for Arabic language. The

f-measure rate from experimental result is 88%. The study presents an effective method for identifying opinion phrases by building Arabic grammatical analyzer with good result and expandable. The future works will be focusing on entering new categories of products and services, support grammatical analyzer with new rules, expand dictionaries, in addition to include other platforms of social media. The sentiment of emoji is one of the future works.

### References

[1] M. N. Al-Kabi, A. H. Gigieh, I. M. Alsamadi, H. A. Washeh and M. M. Haidar, "Opinion Mining and Analysis for Arabic Language," International Journal of Advanced Computer Science and Applications (IJACSA), vol. 5, no. 5, p. 15, 2014.

[2] Dr.S. Murugavalli, U. Bagirathan, R. Saiprassanth and S. Arvindkumar, "Feedback analysis using Sentiment Analysis for E-commerce," International Journal of Latest Engineering Research and Applications (IJLERA), vol. 02, pp. 84-90, 30 March 2017.

[3] F. J. A. P. Mattosinho, "Mining Product Opinions and Reviews on the Web," Master Thesis, W. M. Medieninf, Rer. Nat. Habil, H. C. Alexander Schill (Advisors), Chair of Computer Networks, 2010, July,.

[4] M. El-Masri, N. Altrabsheh, H. Mansour, A. Ramsay, "A web tool for Arabic sentiment analysis," Procedia Computer Science, vol. 117, pp. 38-45, 2017.

[5] B. Pang, L. Lee, "Opinion mining and sentiment analysis," Foundation and Trends in Information Retrieval, Vol. 2 No. 1-2, 1-135, 2008.

[6] S. Alowaidi, M. Saleh, O. Abdulnaja, "Semantic Sentiment Analysis of Arabic Texts," International Journal of Advanced Computer Science and Applications (IJACSA), Vol 8, No. 2, July 2017.

[7] iweb digital advertising agency, "E-commerce stats," iweb123.com, accessed in 2020-10-3 website.

[8] S. Siddiqui, M. Abdul Rahman, S. Daudpota, A. Waqas, "Opinion Mining: Approach to Feature Engineering," International Journal of Advanced Computer Science and Applications (IJACSA), Vol 10, No. 3, 2019.

[9] B. Liu, "Sentiment Analysis and Opinion Mining," Morgan & Claypool, Chicago, 2012.

[10] C. Bhadane, H. Dalal, H. Doshi, "Sentiment analysis: Measuring opinions," Procedia Computer Science, vol. 45, pp. 808-814.

[11] M. Seansuk, P. Songram and P. Chomphuwiset, "Feature-Based Opinion Mining On Smart-Phone Reviews," Proceedings of the 3rd IIAE International Conference on Intelligent Systems and Image Processing, p. 5, 2015.

[12] A. U. R. Khan, M. Khan and M. B. Khan, "Naïve Multi-label classification of YouTube comments using comparative opinion mining," ELSEVIER, Procedia Computer Science 82, vol. 82, pp. 57-64, 12 May 2016.

[13] Weishu Hu, Zhiguo Gong, JingzhiGuo, "Mining Product Features from Online Reviews," IEEE International Conference on E-Business Engineering, 2010.

[14] S. AbdelRahman, M. Elarnaoty, M. Magdy, A. Fahmy, "Integrated Machine Learning Techniques for Arabic Named Entity Recognition," International Journal of Advanced Computer Science and Applications (IJACSA), vol 7, issue 4, No. 3, July 2010.

[15] T. Hassan, M. Ali M, A. Hedar, M. M. Doss, "Mining Social networks' Arabic Slang Comments," Proceedings of IADIS European Conference on Data Mining, 22-24 July.

[16] M. Hammad, M. Al-awaidi, "Sentiment Analysis for Arabic Reviews in Social Network Using Machine Learning," Information Technology, Springer, 2016, pp. 131-139.

[17] A. El-Dine Ali Hamouda, F. El-zaharaa El-taher, "Sentiment Analyzer for Arabic Comments System," International Journal of Advanced Computer Science and Applications (IJACSA), Vol. 4, No. 3, July 2013.

[18] S. Cassab, M. B. Kurdy, "Ontology-based Emotion Detection in Arabic Social Media," International Journal of Engineering Research & Technology (IJERT), Vol. 9 issue 08, August 2020.

[19] A. Daood, I. Salman, N. Ghanem, "Comparison study of automatic classifiers performance in emotion recognition of Arabic social media users," Journal of Theoretical and Applied Information Technology (JATIT), Vol. 95 No. 29, Octobor 2017.

[20] M. Abdul-Mageed, S. Kuebler, M. Diab, "SAMAR: A System for Subjective and Sentiment Analaysis of Social Media Arabic," 3rd Workshop on Computional Approaches to Subjectivity And Sentiment Analysis (WASSA), 2012.

AUTHORS' PROFILE

**Eng. Ghady Alhamad** Born in Syria/Hama July/1987, She obtained bachelor degree in information systems engineering from Syrian Virtual University - Syria December, 2012, Master in Web Science from Syrian Virtual University - Syria Feb, 2013, She started to work as a Software Engineer from 2013. She prefered to work remotely because of her healthy profile with challanging all drawbacks. She worked between UAE and KSA companies.

She is supervisor of graduation projects in Information Systems Engineering Program (Under graduate students) and Master projects (Post graduate) in Web Science Program, Syrian Virtual University since 2015.

She did some special projects: Java compiler, e-commerce sites, web applications, messengers and in augmented reality field View Art in Room. She started in 2020 to publish her works in social media pages on LinkedIn page & Facebook page to leave imprint in this field.

**Ph.D. Mohamad-Bassam Kurdy** Born in Syria/Damascus July/1961, He obtained Master degree in information systems engineering from INPG - France 1986, Ph.D. in Mathematical Morphology from Mines ParisTech -France 1990, He worked at HIAST 1991-2013. He was Head of Computer Sciences at HIAST between 1997 and 2003. Country Manager for Syria of EC project EUMEDIS -Medforist. Actually Professor at SVU, ESC Dijon and ESC Rennes teaching: Advanced Data Mining, Big Data, Information Retrieval, cbIR (content based Image Retrieval) and Supervising many Master student projects (postgraduate).