# A HYBRID METHOD USING LEXICON-BASED APPROACH AND NAIVE BAYES CLASSIFIER FOR ARABIC OPINION QUESTION ANSWERING

**Khalid Khalifa and Nazlia Omar**

Knowledge Technology Group, Center for Artificial Intelligence Technology (CAIT),
Faculty of Information Science and Technology,
University Kebangsaan Malaysia, 43600 Bangi, Selangor, Malaysia

## ABSTRACT

Opinion Question Answering (Opinion QA) is the task of enabling users to explore others opinions toward a particular service of product in order to make decisions. Arabic Opinion QA is more challenging due to its complex morphology compared to other languages and has many varieties dialects. On the other hand, there are insignificant research efforts and resources available that focus on Opinion QA in Arabic. This study aims to address the difficulties of Arabic opinion QA by proposing a hybrid method of lexicon-based approach and classification using Naïve Bayes classifier. The proposed method contains pre-processing phases such as, transformation, normalization and tokenization and exploiting auxiliary information (thesaurus). The lexicon-based approach is executed by replacing some words with its synonyms using the domain dictionary. The classification task is performed by Naïve Bayes classifier to classify the opinions based on the positive or negative sentiment polarity. The proposed method has been evaluated using the common information retrieval metrics i.e., Precision, Recall and F-measure. For comparison, three classifiers have been applied which are Naïve Bayes (NB), Support Vector Machine (SVM) and K-Nearest Neighbor (KNN). The experimental results have demonstrated that NB outperforms SVM and KNN by achieving 91% accuracy.

**Keywords:** Sentiment Analysis, Opinion Question Answering, Naïve Bayes, Lexicon-Based

## 1. INTRODUCTION

With the dramatic expansion of the World Wide Web, the processing of investigating people's opinions has become more accessible and straightforward. Usually, that information would be in a textual mode. Thus, the use of recent technologies such as, web mining and semantic web facilitate the text analyzing which leads to extract the knowledge. Such process called sentiment analysis (Pang and Lee, 2008). Subjectivity is the way that emotions and opinion can be expressed in the language while objectivity refers to the factual phrases (Wiebe *et al*., 1999). The problem of identifying documents whether it is

subjective (yields opinion) or objective (yields fact) is called subjectivity classification. Subjectivity is the way that emotions and opinion can be expressed in the language while objectivity refers to the factual phrases (Wiebe *et al*., 1999). The problem of identifying documents whether it is subjective (yields opinion) or objective (yields fact) is called subjectivity classification. Question answering is the method that provides automated answers for human-generated question. QA is very challenging issue especially when dealing with opinion question answering which has difficulty that lies on the long answer that it need, unlike the factual based questions which have shorter and complete answers (Li *et al*., 2009). The complexity of the

**Corresponding Author:** Khalid Khalifa, Knowledge Technology Group, Center for Artificial Intelligence Technology (CAIT),
Faculty of Information Science and Technology, University Kebangsaan Malaysia, 43600 Bangi, Selangor, Malaysia

Opinion QA resulted from the fact that it combines two very challenging natural languages processing task sentiment analysis and question answering comparing to the traditional QA systems that seek related factual Information to the given question, opinion QA is more complicated due to finding customers sentimental opinion on a specific target.

Arabic is a rich language which has vastly spoken by nearly 300 million people. The Arabic content on the World Wide Web has been increased in recent years. Although it has key challenging, few researches have been done in terms of opinion question answers. There are several reasons that make Arabic opinion QA more challenging. The first reason is the lack of the resources for Arabic opinion QA. Second, it has a very complex morphology as compared to other languages. Moreover, Arabic language has many varieties dialects; this brings another complexity since most writers express their questions and also opinion using local dialect instead of standard Arabic language. So, we end with many written dialects instead of one formal language. This study aims to address the difficulties of Arabic Opinion QA by proposing a hybrid method of lexicon-based approach and classification using Naïve Bayes classifier. The proposed method contains pre-processing phases such as, transformation, normalization and tokenization and exploiting auxiliary information (thesaurus). The experimental results of the proposed method have demonstrated 91% of accuracy. The proposed method has been compared with other existing approaches and shown that the results of this study seem to be promising in the field of Opinion QA.

## 2. RELATED WORK

Devitt and Ahmad (2007) proposed a computable metric of positive or negative orientation in financial news. The authors used cohesion-based text representation algorithm that builds a graph representation of text from part-of-speech POS tagging without disambiguation using Word Net. While (Gigieh et al., 2008) used a customized domain dictionary of Arabic words from social media posts in order to classify the opinion's polarity whether positive or negative. Based on lexicon approach the authors differentiate the formal Arabic language and informal. Abbasi et al. (2008) develop an Entropy Weighted Genetic Algorithms EWGA which is a combination of genetic algorithms and information gain in order to extract feature in movie reviews data set for

English languages and Eastern web forum posts for Arabic language. The results show that EWGA is very competitive approach for feature selection.

This study aims to identify Arabic business reviews that scattered over the internet (approximately 2,000 URLs have been collected). Then a feature selection has been determined which aims to investigate the occurrences number of keywords on each line. Furthermore, they used Ada Boost classifier in order to classify those documents whether it is a review or not. Using a similarity graph which represents similar words (whether in meaning or in polarity) by an edge, the weight of those edges can be an indication to classify the words whether it is positive or negative. The proposed method has achieved highly precision and recall comparing with other mining reviews techniques. Whilst Farra et al. (2010) proposed an approach that can classify both of sentence-level and document-level in terms of Arabic sentiment analysis. In sentence-level the authors used two sub-approaches which are grammatical and semantic. Firstly, they used two features which are sentence type (nominal sentences and verbal sentences) and transition word (words that connect two sentences), secondly (in the semantic side) they used three features which are the frequency of negation words such as don't, didn't and doesn't, the frequency of special characters such as (!) and (?) and the third feature is the frequency of emphasis words such as really, especially and heavily. Eventually, the resulting sentence-level classification will be an input for the document-level in order to classify each sentence. The proposed approach has gained high accuracy that seems to be feasible for both of sentence-level and document-level in terms of sentiment analysis for Arabic language.

El-Halees (2011) proposed a combination of lexicon-based, maximum entropy and k-nearest neighbor with using a corpus of Arabic words in order to classify opinionated Arabic documents. Similarly, Abdul-Mageed and Diab (2011) have built a large corpus for Modern Standard Arabic words from newswire documents annotated on the sentence-level. Their method firstly perform a binary classification in order to classify documents into objective and subjective and then perform a binary classification in terms of classifying subjective documents into positive or negative. Al-Subaihin et al. (2011) this study addresses the obstacles of sentiment analysis when dealing with informal Arabic language. They proposed an unsupervised learning technique which is human-based computing in order to build a lexicon; this would contribute to solve the problem of insufficient informal Arabic corpus. They used a game-based that attracting users to build a lexicon

by tagging words or phrases with its polarity whether positive, negative or neutral. The resulting lexicons will pass to a sentiment analyzer that divide it into sentences, each sentence will be partitioned into set or words and then such words will be matched with its polarity. Elarnaoty *et al*. (2012) proposed a method for extracting the opinion holder in Arabic. The method is based on pattern matching and machine learning. Experimental results on the Arabic Opining Holder corpus show that CRF achieve better results than the pattern matching approach. The authors report 85.52% precision, 39.49% recall and 54.03% F-measure.

# 3. PROPOSED METHOD

The proposed method that has been used consists of several phases. The framework determines the flow of the research methodology phases which include transformation, normalization, tokenization, feature extraction and classification, **Fig. 1** illustrates all those phases.

## 3.1. Transformation

This phase aims to turn the data into an internal representation that enables applying pre-processing and classification steps. Basically, the reviews' comments which are in Arabic language have been collected and stored in text files. Note, that there are several steps have been done in terms of transformation, those steps determines as follows:

### 3.1.1. Turning the Files into UTF-8 Encoding

UTF-8 is an encoding that can be used to simulate any character in the Unicode character. It has been used in order to represent complicated languages such as Arabic language.

### 3.1.2. Arabic Letters

Firstly, all the Arabic letters have to be transformed in order to recognize any Arabic word.

### 3.1.3. Arabic Diacritics

Arabic language has numerous diacritics (tashkil تَشْكِيل), it is vowel marks that indicate how to pronounce letters (e.g., short or long vowel).

### 3.1.4. Definite Articles

It is a set of letters that could be prefix of suffix which indicates to the kind of reference that made by the noun.

### 3.1.5. Arabic Stems

Arabic verbs have numerous forms of derivation. Such forms referred to the origin stem فعل fea'l which means verb.
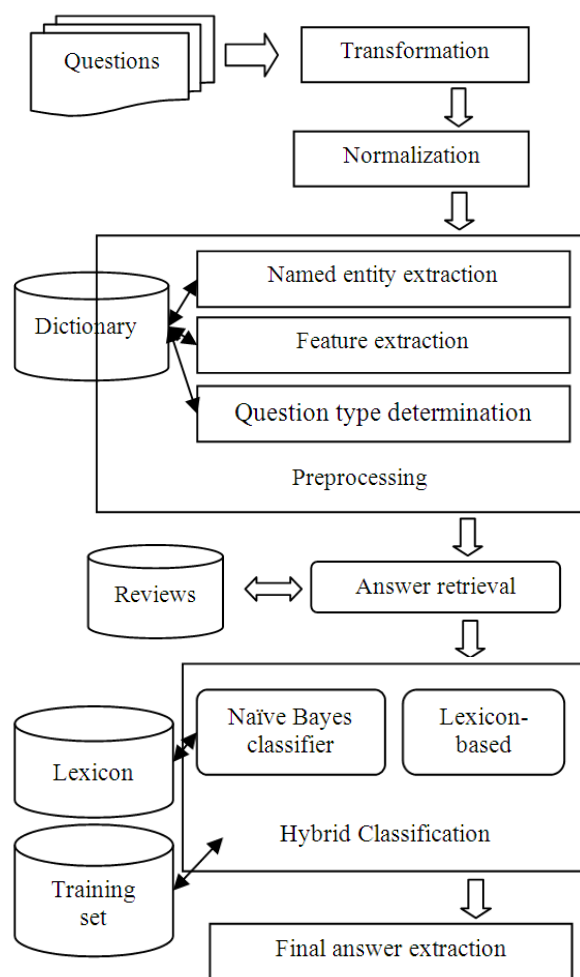


**Fig. 1.** The system design

## 3.2. Pre-Processing

This phase aims to clean, normalize and tokenize the data, in other meaning removing noisy, irrelevant and duplicated data in order to carrying out the classification method on the data. Such data has to be removed due to its effects on the accuracy of the classification. The steps of pre-processing phase are illustrated as follow:

### 3.2.1. Special Character Extraction

Every character seems to be non-Arabic letter has to be removed.

### 3.2.2. Non Arabic Letter Extraction

This process aims to remove special character (e.g., _+-:;""?).

### 3.2.3. Digit Extraction

Digits from 0 to 9 have to be removed.

### 3.2.4. Diacritics Extraction

As mention before, Arabic language has several forms of diacritics thus; those diacritics have to be removed.

### 3.2.5. Definite Articles Extraction

As mention before, those articles could be prefix of suffix (**Table 3**). Such articles have to be removed.

### 3.2.6. Tokenization

Tokenizing the data indicates to deleting multiple spaces.

### 3.2.7. Stemming

As mention before, Arabic language has numerous forms of verbs thus; those derivations have to be stemmed to its own roots.

### 3.2.8. Stop-Words Extraction

Stop words are words that connect the sentences which seem to be irrelevant and have to be removed.

### 3.2.9. Filtering Arabic Letters

There are some letters in Arabic that have several forms, so that those letters have to be unified.

### 3.3. Question Analysis

In this section, several steps of analysis have been conducted in order to acquire more informative retrieval of the questions. Each question yields one or more indication for product, service or an industry. Therefore, utilizing those features may facilitate to retrieve more relevant questions. Those analytical steps determined as follows:

### 3.3.1. Named Entity Extraction

In this step, it aims to identify the name of the entities that appears in the questions for instance, questions that ask about a particular industry such as hotels or resorts have to be analyzed in order to extract those entities. This can facilitate the analysis of the question. **Table 1** illustrates an example of this step.

### 3.3.2. Feature Extraction

In this step, feature extraction aims to identify the features that appear in the questions for instance, questions that ask about a specific service in a particular industry such as, room service which indeed located in hotels, have to be analyzed in order to extract those features (room service). **Table 2** illustrates an example of extracting the feature type.

### 3.3.3. Question Type Determination

This step attempts to identify the type of the question. There are several forms of questions such as descriptive or comparative. **Table 3** shows an example.

### 3.4. Lexicon-Based

Sentiment lexicon or so called senti-lexicons is the process where each word is associated with a polarity score which indicates the orientation of the word (positive or negative). The sentiment lexicon is the most sensitive resource for most sentiment analysis algorithms (Feldman, 2013). In the following, an illustration for designing Arabic sentiment lexicon has been described in details. **Table 4** shows the structure of the proposed Arabic sentiment lexicon.

As seen in **Table 4**, syntactic refers to the annotated part-of-speech for each words (adjectives, adverb, nouns and verbs). The common syntactic that usually used is the adjectives. Whereas, score is refer to the degree of polarity from most bad to most good which has been ranged between -5 to 5. Eventually, inflections forms refer to the forms that the word can be formulated whether for singular, female, dual, or plural.

### 3.5. Classifier Selection

This study aims to take the advantage of Naïve Bayes classifier in order to combine it with the lexicon based. Naïve Bayes has several advantages but the main one is the ability to classify rapidly without using a large training set. This algorithm is very effective in terms of classifying documents. It has been vastly proposed in literature such as (Melville *et al*., 2009; Pang and Lee, 2004; Tan *et al*., 2009). It uses the statistics in order to perform probabilities classification. Basically, it aims to analyze the absence and presence of a particular feature in order to independently classify feature using probabilities. It is very effective when treating words that have probabilities to be opinion or not such as, adjectives or adverbs.

**Table 1.** Sample of named entity extraction

| Question | Named entity |
|---|---|
| ماهو تقييمك لفندق راديسون بلو باي ريزورت؟ | فندق راديسون بلو |
| What is your feedback about Radisson Blue hotel? | Radisson Blue hotel |
| ما هو تقييمك لفندق موفينبيك؟ | فندق موفينبيك |
| What is your feedback about MovenPick hotel? | MovenPick hotel |

**Table 2.** Feature extraction

| Question | Feature type |
|---|---|
| ما هو مستوى الغرف في فندق راديسون بلو باي ريزورت؟ | مستوى الغرف |
| What do you think about the room service in Radisson Blue Hotel? | Room service |
| ما هو تقييمك للمطعم في فندق كيمبنسكي؟ | مطعم فندق كيمبنسكي |
| What is your feedback about the restaurant of Kempinski Hotel? | Kempinski hotel's restaurant |

**Table 3.** Question type determination

| Question | Question type |
|---|---|
| برأيك هل مطعم فندق موفينبيك أفضل من مطعم فندق راديسون بلو؟ | Comparative question |
| Do you think that MovenPick's hotel's restaurant better than Radisson Blue hotel's restaurant? | |
| ما هو رأيك في حمام السباحة بفندق راديسون بلو؟ | Descriptive question |
| What is your feedback about the swimming pool on Radisson Blue Hotel? | |

**Table 4.** Structure of Arabic lexicon

| Word | Syntactic | Synonyms | Score | Inflections forms | Meaning | Polarity |
|---|---|---|---|---|---|---|
| جيد | Adjective | منيح ؛ تمام ؛ زين | 4 | جيدون ، جيدان ، جيدات | Good | Positive |
| جميل | Adjective | حلو ؛ رائع ؛ رهيب | 4 | جميلون ، جميلان ، جميلات | Beautiful | Positive |
| مقرف | Adjective | وحش ؛ مو زين ؛ سيء | -4 | مقرفون ، مقرفان ، مقرفات | Disgusting | Negative |

It can be formulated as follows Equation 1:

$$P(C_i \mid d) = \frac{P(C_i)P(d \mid C_i)}{P(d)} \qquad (1)$$

where, $P(C_i|d)$ is the posterior probability of class Ci given a new document d, $p(C_i)$ is the probability of class Ci which can be calculated by Equation 2:

$$P(C_i) = \frac{N_i}{N} \qquad (2)$$

where, Ni is the number of documents that associated with class Ci and N is the number of classes, $P(C_i|d)$ is the probability of a document d given a class Ci and P(d) is the probability of document d. Because the independence assumption of NB, The probability of document d can be calculated by Equation 3:

$$P(C_i \mid d) = P(C_i)\prod_{k=1}^{n}(t_k \mid C_i) \qquad (3)$$

where, tk is a feature that co-occurs with class Ci and also we can calculate $(t_k|C_i)$ by Equation 4:

$$P(t_k \mid C_i) = \frac{1 + nk_i}{l + \sum_{h=1}^{l} n_{hi}} \qquad (4)$$

where, n hi is the total number of documents that contain feature tk and belong to class Ci, l is the total number of distinct features in all training documents that belong to class Ci. NB calculates posterior probability for each class and then assigns document d to highest posterior probability's class, i.e., Equation 5:

$$C(d) = \underset{i=}{\arg\max} \mid C \mid (P(C_i \mid d)) \qquad (5)$$

The reason that this study used NB as classifier is that NB has the ability to classify objects based on its features independently, in other meaning NB focuses on some features regardless of the absence or presence of others for instance, a negative answer may consist of negation, adjective or adverb NB considers each of these features to

contribute independently to the probability that this answer is negative regardless of the presence or absence of the other features. This scenario is very common in the Arabic where all features may not exist. Therefore, by applying NB on the reviews comments in our study, the classification has been done by identifying the number of positive and negative reviews. Note that, the opinions presented with its associated feature, location and all the specific details. **Figure 2** shows the algorithm of the proposed method.
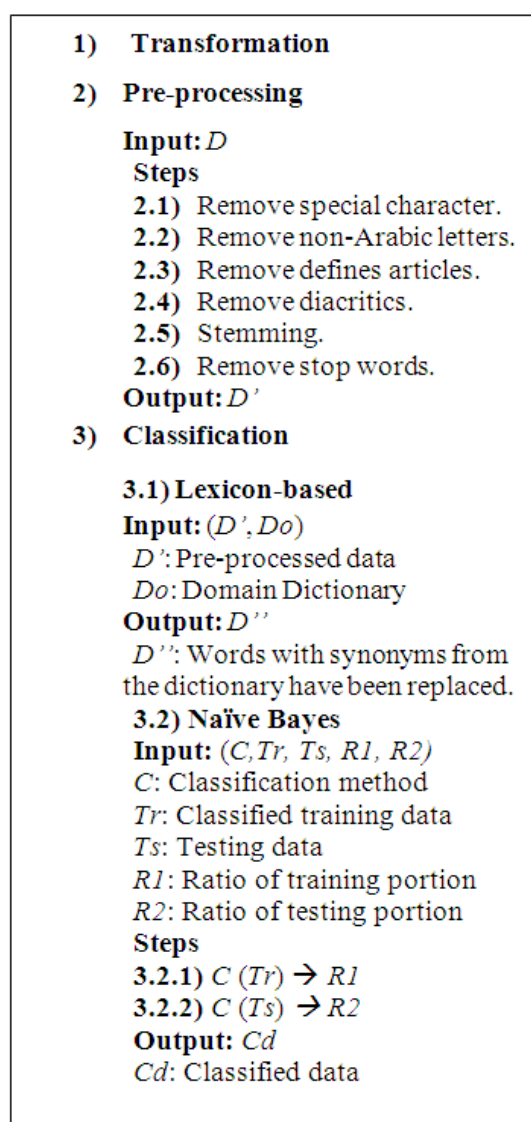
```
1)  Transformation

2)  Pre-processing
    Input: D
    Steps
    2.1)  Remove special character.
    2.2)  Remove non-Arabic letters.
    2.3)  Remove defines articles.
    2.4)  Remove diacritics.
    2.5)  Stemming.
    2.6)  Remove stop words.
    Output: D'

3)  Classification

    3.1) Lexicon-based
    Input: (D', Do)
     D': Pre-processed data
     Do: Domain Dictionary
    Output: D''
     D'': Words with synonyms from
     the dictionary have been replaced.
    3.2) Naïve Bayes
    Input: (C, Tr, Ts, R1, R2)
    C: Classification method
    Tr: Classified training data
    Ts: Testing data
    R1: Ratio of training portion
    R2: Ratio of testing portion
    Steps
    3.2.1) C (Tr) → R1
    3.2.2) C (Ts) → R2
    Output: Cd
    Cd: Classified data
```

**Fig. 2.** Algorithm of the proposed method

## 4. RESULTS

Basically, the results consist of two parts which are lexicon-based experiment results and classifiers experiments results. In the first experiment, several types of features and investigating their effect on performance of the Arabic opinion question answering have been taken place. The aim is to efficiently integrate those features with the most accurate classification algorithm. There are several features have been used on the lexicon-based approach phase. Firstly, sentiment words polarity have been tested as a feature set (F1). Then, named entity such as, hotel and resort have been also tested as a feature set (F2-F3). Whereas, the feature set (F4-F7) refers to the special features that indicate type of services such as, room service. Eventually, the feature set (F8-F9) addresses the determination of the question whether it is comparative or descriptive. **Table 5** illustrates those feature sets in details.

The second part of results is the classifiers results. Basically, three classifiers have been carried out with the lexicon-based approach. Such classifiers are Naïve Bayes NB, Support Vector Machine SVM and K-nearest neighbor KNN. The evaluation has been performed using the macro-average F-measure. Note that, the three classifiers have been performed firstly without the lexicon-based. **Table 6** shows the results of the classifiers without lexicon-based approach.

Sequentially, each classifier has been combined with the lexicon-based approach independently in order to investigate the role of lexicon-based approach. **Table 7 to 9** show the results of those classifiers with lexicon-based approach.

## 5. DISCUSSION

As shown in **Table 7-10**, there is a remarkable enhancement when using lexicon-based approach which demonstrates the significant impact of utilizing lexicon-based. On other hand, **Table 10** shows that NB has outperformed SVM and KNN by achieving 0.91 F1-measure for positive, 0.91 F1-measure for negative and 0.91 for the macro F-measure. As a result, NB has been selected to be the best classifier for the combination with lexicon-based approach. Therefore, the proposed hybrid method contains NB classifier and lexicon-based approach. The proposed method has demonstrated a remarkable increase in the Macro-averaging of F-measure. However, since this study uses a different dataset compared to other previous work, the results cannot be compared directly with their results.

**Table 5.** Feature sets of lexicon-based

| Feature set name | Feature | Feature name |
|---|---|---|
| Sentiment words polarity | F1 | The weighted polarity of words |
| Named entity | F2 | The proportion of opinions that indicate hotel |
| | F3 | The proportion of opinions that indicate resort |
| Special features | F4 | The proportion of opinions that indicate room services |
| | F5 | The proportion of opinions that indicate restaurant services |
| | F6 | The proportion of opinions that indicate general services |
| | F7 | The proportion of opinions that indicate quality services |
| Question determination | F8 | The proportion of comparative opinions |
| | F9 | The proportion of descriptive opinions |

**Table 6.** Results of three classifiers without lexicon-based approach

| | Macro F-measure | F1 Measure for sentiment positive | F1 Measure for sentiment negative |
|---|---|---|---|
| NB | 0.82 | 0.85 | 0.80 |
| SVM | 0.80 | 0.80 | 0.81 |
| KNN | 0.80 | 0.79 | 0.82 |

**Table 7.** Results of NB with lexicon-based approach

| | NB | | |
|---|---|---|---|
| | Macro F-measure | F1 Measure for sentiment positive | F1 Measure for sentiment negative |
| F1 | 0.92 | 0.92 | 0.93 |
| F2-F3 | 0.89 | 0.89 | 0.90 |
| F4-F7 | 0.92 | 0.92 | 0.93 |
| F8-F9 | 0.91 | 0.91 | 0.91 |
| Macro-averaging | 0.91 | 0.91 | 0.91 |

**Table 8.** Results of SVM with lexicon-based approach

| | SVM | | |
|---|---|---|---|
| | Macro F-measure | F1 Measure for sentiment positive | F1 Measure for sentiment negative |
| F1 | 0.86 | 0.86 | 0.87 |
| F2-F3 | 0.88 | 0.88 | 0.88 |
| F4-F7 | 0.88 | 0.89 | 0.87 |
| F8-F9 | 0.88 | 0.90 | 0.87 |
| Macro-averaging | 0.87 | 0.88 | 0.87 |

**Table 9.** Results of KNN with lexicon-based approach

| | KNN | | |
|---|---|---|---|
| | Macro F-measure | F1 Measure for sentiment positive | F1 Measure for sentiment negative |
| F1 | 0.84 | 0.84 | 0.85 |
| F2-F3 | 0.85 | 0.87 | 0.83 |
| F4-F7 | 0.87 | 0.86 | 0.89 |
| F8-F9 | 0.87 | 0.88 | 0.86 |
| Macro-averaging | 0.85 | 0.86 | 0.85 |

**Table 10.** Macro-average results for the three classifiers with lexicon-based approach

| | Macro-average precision | F1-measure for positive | F1-measure for negative |
|---|---|---|---|
| NB | 0.91 | 0.91 | 0.91 |
| SVM | 0.87 | 0.88 | 0.87 |
| KNN | 0.85 | 0.86 | 0.85 |

# 6. CONCLUSION

This study proposed a hybrid classification method consist of Naïve Bayes classifier and lexicon-based approach in order to build Arabic Opinion question answering. The dataset has been collected from the comments of Jordan hotels and resorts' residents. After

transformation and preprocessing phases, three classifiers have been carried out with the lexicon approach. Those classifiers are NB, SVM and KNN. Eventually, NB has demonstrated best results of macro-average F-measure of 91%.Therefore, NB was selected to be combined with the lexicon-based approach in order to solve the problem of sentence-level sentiment analysis.

# 7. REFERENCES

Abbasi, A., H. Chen and A. Salem, 2008. Sentiment analysis in multiple languages: Feature selection for opinion classification in Web forums. ACM Trans. Inf. Syst., 26: 1-34. DOI: 10.1145/1361684.1361685

Abdul-Mageed, M. and M.T. Diab, 2011. Subjectivity and sentiment annotation of modern standard Arabic newswire. Proceedings of the 5th Linguistic Annotation Workshop, (LAW' 11), pp: 110-118.

Al-Subaihin, A.A., H.S. Al-Khalifa and A.S. Al-Salman, 2011. A proposed sentiment analysis tool for modern Arabic using human-based computing. Proceedings of the 13th International Conference on Information Integration and Web-based Applications and Services, Dec. 05-08, Hue City, Viet Nam, pp: 543-546. DOI: 10.1145/2095536.2095651

Devitt, A. and K. Ahmad, 2007. Sentiment polarity identification in financial news: A cohesion-based approach. Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics, (ACL' 07), Czech Republic, pp: 984-991.

El-Halees, A., 2011. Arabic opinion mining using combined classification approach. Proceeding of the International Arab Conference on Information Technology, (CIT' 11), Azrqa, Jordan, pp: 1-8.

Elarnaoty, M., S. AbdelRahman and A. Fahmy, 2012. A machine learning approach for opinion holder extraction in Arabic language.

Farra, N., E. Challita, R.A. Assi and H. Hajj, 2010. Sentence-level and document-level sentiment mining for Arabic texts. Proceedings of the IEEE International Conference on Data Mining Workshops, Dec. 13-13, IEEE Xplore Press, Sydney, NSW, pp: 1114-1119. DOI: 10.1109/ICDMW.2010.95

Feldman, R., 2013. Techniques and applications for sentiment analysis. Commun. ACM, 56: 82-89. DOI: 10.1145/2436256.2436274

Gigieh, A., M. Al-Kabi, I. Alsmadi, H. Wahsheh and M. Haidar, 2008. Building and evaluating an opinion analysis tool for standard and colloquial Arabic language. Yarmouk University.

Li, F., Y. Tang, M. Huang and X. Zhu, 2009. Answering opinion questions with random walks on graphs. Proceedings of the Joint Conference of the 47th Annual Meeting of the Association for Computational Linguistics, (ACL' 09), Stroudsburg, PA, USA, pp: 737-745.

Melville, P., W. Gryc and R.D. Lawrence, 2009. Sentiment analysis of blogs by combining lexical knowledge with text classification. Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Jun. 28-Jul. 01, Paris, France, pp: 1275-1284. DOI: 10.1145/1557019.1557156

Pang, B. and L. Lee, 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics, (ACL' 04), Stroudsburg, PA, USA, pp: 271-271. DOI: 10.3115/1218955.1218990

Pang, B. and L. Lee, 2008. Opinion mining and sentiment analysis. Found. Trends Inform. Retr., 2: 1-135. DOI: 10.1561/1500000011

Tan, S., X. Cheng, Y. Wang and H. Xu, 2009. Adapting Naive Bayes to Domain Adaptation for Sentiment Analysis. In: Advances in Information Retrieval, Boughanem, M., C. Berrut, J. Mothe and C. Soule-Dupuy, (Eds.), Springer Berlin Heidelberg, pp: 337-349. DOI: 10.1007/978-3-642-00958-7_31

Wiebe, J.M., R.F. Bruce and T.P. O'Hara, 1999. Development and use of a gold-standard data set for subjectivity classifications. Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics on Computational Linguistics, (LCL' 99), Stroudsburg, PA, USA, pp: 246-253. DOI: 10.3115/1034678.1034721