

# Discovering Semantic and Sentiment Correlations using Short Informal Arabic Language Text

Salihah AlOtaibi

Information Systems Department, College of Computer and  
Information Sciences  
Al Imam Mohammad Ibn Saud Islamic University (IMSIU)  
Riyadh, KSA

Muhammad Badruddin Khan

Information Systems Department, College of Computer and  
Information Sciences  
Al Imam Mohammad Ibn Saud Islamic University (IMSIU)  
Riyadh, KSA

**Abstract**—Semantic and Sentiment analysis have received a great deal of attention over the last few years due to the important role they play in many different fields, including marketing, education, and politics. Social media has given tremendous opportunities for researchers to collect huge amount of data as input for their semantic and sentiment analysis. Using twitter API, we collected around 4.5 million Arabic tweets and used them to propose a novel automatic unsupervised approach to capture patterns of words and sentences of similar contextual semantics and sentiment in informal Arabic language at word and sentence levels. We used Language Modeling (LM) model which is statistical model that can estimate the distribution of natural language in effective way. The results of experiments of proposed model showed better performance than classic bigram and latent semantic analysis (LSA) model in most of cases at word level. In order to handle the big data, we used different text processing techniques followed by removal of the unique words based on their role Informal Arabic, Big Data, Sentiment analysis, Opinion Mining (OM), semantic analysis, bigram model, LSA model, Twitter vance to problem.

**Keywords**—Opinion Mining; Sentiment analysis; semantic analysis; Twitter; Informal Arabic

## I. INTRODUCTION

The last decade has seen a huge increase in the number of internet users in Middle East. This growth has helped in enriching the amount of Arabic content on website. There are wide numbers of users that use the social networks. They use social media in order to share various kinds of resources, express their opinions, thoughts, and messages in real time [1]. Since most of users use informal Arabic in the world of social media, the task of semantic and sentiment analysis becomes more sophisticated. Different Arabic Dialects are another challenge [2]. One of the main challenges is the limited number of researches that focus on the informal Arabic sentiments analysis. This motivated us to focus on the problems that exist in the realm of informal Arabic semantic and sentiment analysis thus encouraging the researchers to participate more in this field. Sentiment analysis, also called opinion mining, is the field of study that extracts and analyzes people's opinions about products, services, individuals, event, issues, to name but a few categories [3][4][5]. An opinion can be a positive or negative or neutral sentiment, attitude, emotion, or appraisal. For small corpus of data, it is possible to use humans for annotation but for big data, the formulation of training and test data is very expensive and almost

impossible. Although a tweet is small piece of data but to annotate them when they are millions followed by application of machine learning techniques and then analyzing classification models to understand the polarity of different words is pretty difficult and expensive job.

This work is neither related to supervised learning nor it use existing semantic resources like Arabic WordNet due to informal nature of Arabic text in tweets. The proposed approach does not depend on the syntactic structure of tweets, it extracts patterns from the contextual semantic and sentiment similarities between words in a given tweet corpus. Contextual semantics are based on the proposition that meaning can be extracted from words co-occurrences [6]. The LM model gives a probability distribution—or  $P(s)$ —over sequences of words ( $w_i$ ). The goal of LM is to build a statistical model that can estimate the distribution of natural language in effective way [7] [8]. It uses a number of types of matrices, such as the unigram, bigram, and trigram. The bigram matrix is sometimes referred to as the word co-occurrence matrix [9][10]. In this study, we use a bigram matrix method for document representation. In the bigram matrix, each row represents a word ( $w_i$ ), and each column represents the first preceding word ( $w_j$ ) of  $w_i$  where  $j = i-1$ . Each cell gives the co-occurrence frequency ( $a_{ij}$ ) of the word sequence  $w_j w_i$  in the corpus [9] [11][12].

The paper is organized in few sections to describe further details of the work to extract semantic and sentiments from the huge corpus of Arabic tweets. Section 2 outlines the related work done in this area. In section 3, describes the methodology of the work. In Section 4, discuss the experiments results. Finally, in the brief Section 5 we will make concluding remarks.

## II. RELATED WORK

This section provides a literature review for the field of sentiment and semantic analysis, focusing mainly on informal Arabic language.

### A. Sentiment analysis in informal Arabic language

Duwairi, Marji, Sha'ban, and Rushaidat look at Arabic dialects, Arabism and emoticons. At the normalization stage, they add new step, which is to convert Arabic dialect to Modern Standard Arabic (MSA) by mapping dialect words on to MSA stems. Their study applied three different classifiers: Support Vector Machines (SVM), Naive Bayes (NB), and

KNN. The accuracy of the SVM was found to be slightly lower than that of NB [13]. Both [14] and [2] have produced applications for Arabic sentiment analysis in order to classify Arabic tweets. They used the SVM and Naive Bayes classifiers, and also tried to classifiers together. Itani, Hamandi, Zantout, and Elkabani have studied the use of informal Arabic on Facebook. Their corpus contained eight different dialects; namely, Lebanese, Egyptian, Syrian, Iraqi, Libyan, Algerian, Tunisian, Sudanese, and Saudi. They built a classifier model using the Naive Bayes classifier. Accuracy was measured by comparing human and automatic classification results [15]. Other researches have focused upon the lexicon-based approach, which, typically, is used less often in Arabic sentiment analysis because of the low number of existing Arabic sentiment lexicons. The main challenge here is in building lexicons for informal words, as [1] [16] [17] and [18]. These studies encourage researchers to contribute more extensively to the field. El-Beltagy and Ali (2013) use the semantic orientation approach (SO) to determine Arabic Egyptian polarities, using two data sets: a Twitter data set and a Comments data set. The experiment showed that SO is effective, especially within the context of Twitter [16]. One of the latest sentiment analysis studies has been conducted by [17]. They analyzed three constructed lexicons, one manual and two automatic, designing a lexicon-based model for sentiment analysis. The result of performance is 74.6% is very encouraging. However, some interesting research has been undertaken that uses semantic analysis methods with the aim of improving the sentiment model. Unfortunately, in terms of our context, these studies focus on the English language. Saif, He, and Alani demonstrate the importance of using semantic features when determining the positive or negative sentiments in tweets. In their study, they used both tweet- and entity-level sentiment analysis [19]. They also propose a further study capturing the patterns of word with similar contextual semantics and sentiments in tweets [6]. [20] used a vector space model that learns word representations in order to capture semantic and sentiment words.

### B. Semantic analysis in informal Arabic language

This section offers an overview of some studies that have applied semantic analysis to Arabic language data sets. The amount of Arabic language documents available online is increasing with time. It is difficult for researchers to handle huge volumes of relevant texts documents. For this reason, Arabic document clustering is an important task for achieving perfect outcomes with Information Retrieval (IR) programs, thus satisfying a researcher's needs. Froud, Lachkar, and Ouatic have proposed a method for improving document categorization by using the topic map method, based on a method similar to document clustering. Their method was found to be quite effective for clustering documents, when compared with evaluation methods involving human beings [21]. Other study has sought to group the semantic features of Arabic web pages, clustering them based on their similarities,

with the help of the Arabic VerbNet lexicon. The researchers collected a corpus from the archives of digital Arabic newspapers [22]. Other researchers propose the use of an Arabic language model for speech recognition and machine translation tasks [23] [24]. Notably, Sarikaya et al. introduced the joint morphological-lexical language model (JMLLM), which takes advantage of Arabic morphology, being designed for inflected languages in general (and Arabic in particular). They have used their system to conduct experiments into dialectal Arabic (Iraqi Arabic specifically). The results showed that JMLLM offers encouraging improvements when compared with base-line word- and morpheme-based trigram language models [23]. Latent semantic analysis (LSA) is promoted by many researchers, such as Froud, Lachkar, and Ouatic (2012), who offer an LSA method that uses a variety of distance functions and similarity measures to determine the similarities between Arabic words. Their study compares the results for the use of the model with and without stemming. It was found that stemming affects the obtained results negatively when using the LSA model [25]. The same authors also used their system to produce new results for their previous experiment by comparing stemming and light stemming. The results showed that the light stemming approach out-performed the stemming approach because the latter affects the meanings of words [26]. In the medical domain, the LSA method has been used to predict protein-protein interaction, based on the Arabic semantic analysis model. This method was used to help the researchers understand how and why two proteins interact because protein pairs may interact if they contain similar or related Arabic words. This new method was compared with two other successful methods – namely, PPI-PS and PIPE, and higher accuracy was achieved with the new methods. This research gives insight, therefore, into the importance of semantic analysis, as this method achieved more accurate results than other successful methods [27].

### III. METHODOLOGY

A novel approach to improve the performance measures of informal Arabic language sentiment analysis is proposed to analyze the semantics and sentiment of user-generated text at the word and sentence level. We automatically capture patterns of words of similar contextual semantics and sentiment in tweets. The proposed approach does not depend on the syntactic structure of tweets; instead, it extracts patterns from the contextual semantic and sentiment similarities among words in a given tweet corpus. Contextual semantics are based on the proposition that meaning can be extracted from words' co-occurrences. Evaluating the proposed approach by comparing the results of the proposed approach with the results of the classic bigram and LSA approach. Figure 1 illustrates the semantic sentiment analysis model for informal Arabic. An overview of the framework's four stages, as depicted in Figure 1, is presented in this section. The four stages of proposed framework are as follows:

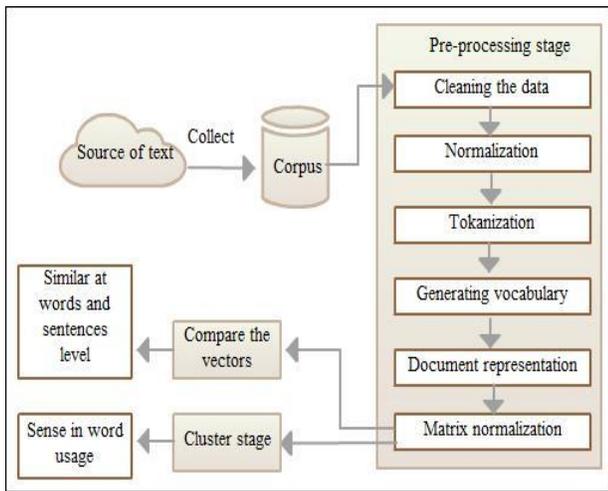


Fig. 1. Framework for unsupervised clustering methodology

### A. Dataset collection

As seen in Figure 1 the first step is document gathering, which is conducted in order to build a corpus. We had to collect our own specialized data (i.e., tweets generated in KSA). For this research the Twitter microblog is one of the best resources for collecting dataset. To collect the Arabic tweets, we used Twitter's stream API in order to avoid the problems of bias and excessive time consumption that can occur when collecting the data manually. The corpus contained 4,425,003 tweets that was saved in a database.

The data collection began on July 7, 2014. The duration of the data collection period coincided with the following events: the month of Ramadan, the FIFA World Cup, and Eid al-Fitr.

### B. Pre-processing

The preprocessing stage is very important in achieving good results from text mining. In context of big data, it can also be seen as a preventive measure to handle the curse of dimensionality. Thus we created our own text preprocessing scheme to deal with informal Arabic language (i.e., Saudi dialect). The text preprocessing stage contains the following four steps:

#### 1) Cleaning the dataset

The cleaning process is used to remove all of the following cases:

- separate any non-Arabic word followed by Arabic word by single space, for example,  
Noor هلا -> Noor هلا
- separate any Arabic word followed by non-Arabic word by single space, for example,  
email ارسل -> email ارسل
- replace all URLs with the symbol URL
- replace all emails with the symbol EMAIL
- replace all time formats with the symbol TIME

- replace all date formats with the symbol DATE
- replace all numbers with the symbol NUMBER
- remove repeated characters, for example, nooo-> noo
- remove repeated sequences no no no no no -> no no
- Separate symbol sequences, for example, ?!! -> ? ! !

We used this process for cleaning in order to reduce the corpus size and noise, while also ensuring that context of the tweet remains unchanged.

#### 2) Normalization

The normalization process is manipulating the text to produce consistent form, by converting all the various forms of a word to a common form. Table.1 shows the all normalization cases that we handled in study experiment.

TABLE I. NORMALIZATION CASES

| Rule     | Example            |
|----------|--------------------|
| Tashkeel | المؤمنين->المؤمنين |
| Tatweel  | الله->الله         |
| Alef     | اor ا!->ا          |
| Heh      | هor ه->ه           |

### C. Tokenization

The tokenization process was performed for each tweet in order to divide the tweet into multiple tokens based on whitespace characters. The corpus was divided into 1,383,012 unique words.

### D. Generating vocabulary

This process was used to build a list of vocabulary words that used the list of pairs (i.e., the word and its counts); the word order was arranged alphabetically. This resulted in 1,383,012 unique words. Then, to avoid out of memory problem, we reduced the vocabulary size to 13,696 words by deleting the words that appeared fewer than 400 times in the corpus, which equals 84% of the corpus. The computational and storage resources largely determined the frequency limit.

### E. Document representation

In this step, these numerical data were transformed into vectors. The Bigram matrix was used to implement this task. The bigram matrix only contains numerical data. The Bigram matrix denoted by  $X_{13696 \times 13696}$  has size of  $13696 \times 13696$ . Each entry in the matrix represents the frequency (i.e., how many times  $w_j$  came before  $w_i$  in the corpus). Figure 2 illustrates the process. The matrix contains the co-occurrence frequency for the words before and after; if we take sequence  $w_2 w_1$ , then word  $w_2$  came before  $w_1$ , and if we take sequence  $w_1 w_2$ , then word  $w_1$  came before  $w_2$ . In other words,  $w_2$  came after  $w_1$ .

While the matrix is square, if we take the transpose of  $X$  (i.e.,  $X^T$ ) we will be able to determine how many times  $w_j$  came after  $w_i$ . Then, concatenate the two matrixes together to make a new matrix and to make each vector contain the before and after frequency value. The new matrix is  $X = [XX^T]$ . The new size is  $n \times 2m$ , where  $n = 13696$  and  $m = 13696$ .

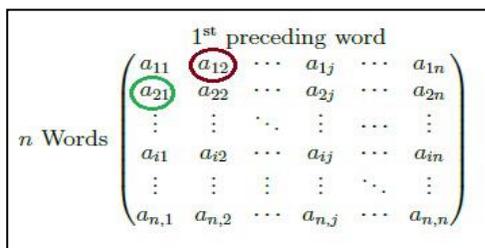


Fig. 2. Bigram matrix

F. Normalization

The normalization helps prevent attributes with large ranges from outweighing attributes with smaller ranges (Jonker, Petkovic, Lin, & Selcuk Candan, 2004). The bigram matrix for any given training corpus is sparse; most of the elements in the matrix are zero. This task of re-assessing some of the zero value and assigning them non-zero values is called smoothing. Then add one to all the counts in the matrix called X. This algorithm is called add-one smoothing (Jurafsky & James, 2000).

Then, used the column wise method to normalize the columns in matrix X by summing the elements in each column, i.e.  $\sum_{i=1}^n a_{ij}$ , where j is the column number. Then divide each element in the matrix with the perspective column sum, i.e.  $a_{ij} / \sum_{i=1}^n a_{ij}$ . (Novak & Mammone, 2001) Then, the based 2 log is calculated for all elements in the matrix X, to make the data more normally distributed (Zhai, 2008). Then, the z-score for all elements X is computed by subtracting the mean and dividing by the standard deviation. This should first by apply in the columns level and then the rows level.

$$x_{norm} = \frac{x - \mu}{\sigma}$$

Where x is referring to the score,  $\mu$  refers to mean and  $\sigma$  refers to standard deviation.

G. Clustering Stage

After normalizing the numerical data, the words dimensions were reduced by applied K-means algorithm to categorize the words by setting k = 200. After many experiments we arrived at k=200 as the best result for word clustering to capture patterns of words of similar contextual semantics and sentiment in tweets (see section 4 to see the experiment result). Figure 3 illustrates part of bigram matrix after normalization

|    | A       | B       | C       | D       | E       | F       | G       |
|----|---------|---------|---------|---------|---------|---------|---------|
| 1  | -2499.5 | -2422.1 | -777.85 | -2387.2 | -2324.3 | -4497.8 | 2174.9  |
| 2  | -2929.8 | -2968.2 | -1123.1 | -3114.6 | -2736.6 | -3829.5 | 1603.2  |
| 3  | -2671.7 | -2496   | -148.86 | -2671.9 | -2359.9 | -2910.8 | 2380.4  |
| 4  | -926.15 | -928.25 | -964.55 | -827.53 | -406.08 | -1382.5 | -538.74 |
| 5  | 164.68  | 17.282  | 1158.1  | -282.03 | 1016.8  | 414.08  | 595.85  |
| 6  | -286.54 | -447.94 | -23.38  | -627.81 | -46.362 | -98.742 | -208.17 |
| 7  | -351.95 | -381.5  | -9.7788 | -472.51 | 327.51  | -478.33 | -98.13  |
| 8  | -1392.5 | -1417.1 | -1676.2 | -1273   | -741.05 | -2205.5 | -462.82 |
| 9  | -434.74 | -580.88 | -182.7  | -566.88 | 210.84  | -618.31 | 38.445  |
| 10 | 198.17  | -21.362 | 811.79  | -229.82 | 641.53  | 98.06   | 562.57  |
| 11 | -582.16 | -523.38 | 373.11  | -743.49 | -4.1742 | -632.66 | 629.31  |
| 12 | -14.064 | -363.34 | 1537.8  | -406.97 | 214.62  | -336.83 | 1501.8  |
| 13 | 429.52  | 234.68  | 1233.8  | 170.34  | 1031.3  | 505.72  | 403.42  |
| 14 | -1540.7 | -1265.5 | -1625.9 | -1305.6 | -1293.7 | -1182.3 | -1045.2 |

Fig. 3. Screenshot of bigram matrix after normalization

To find the similar contextual semantics and sentiment for the sentence level, we calculated the average of the words' vectors that appeared in the sentence in order to get a new vector for the sentence. If we have sentence  $S_i = \{w_1, w_2, \dots, w_n\}$ .

$$\text{Sentence vector} = \frac{\sum_{i=1}^n V_{w_i}}{\text{Total number of words}}$$

Where  $V_{w_i}$  denote to the value of the words' vector in the sentence  $S_i$ . For example, if we have a tweet: "I love Mac products". The vector of each word:

"I": [1,0,3], "love": [1,1,5], "Mac": [2, 0,3], "products": [0, 0, 2]

The sum = [4, 1, 13]

Sentence vector = [4,1, 13]/ 4 = [1, 0.20, 3.25]

One challenge in clustering of short text (e.g., tweets) is that exact keyword matching may not work well (Aggarwal & Zhai, 2012). This research overcomes this challenge and extracts patterns automatically of words of similar contextual semantics and sentiment in tweets.

H. The Model Validation

This stage evaluates the model by comparing the model results with the results of the bigram model and LSA model. The bigram model used the same vocabulary size to build the matrix and also used the same normalization process. The bigram matrix denoted as matrix V with size 13696x13696. The LSA model used feature extraction TF-IDF, and set the SVD rank to used feature extraction TF-IDF, and set the SVD (singular value decomposition) rank to  $K = 100$ . We try to set the  $K = 200$ , but hardware limitation did not permit to perform experiments with this setting.

IV. EXPERIMENTS RESULTS

In this section, we will present necessary information about study experiments along with the results. The tailored bigram model  $[XX^T]$  that was discussed in previous section, was used for experiments. The words dimensions were reduced by utilizing K-means clustering to analyze the semantics and sentiments of user-generated text at word and sentence levels. The proposed method was then compared with the bigram [X] and LSA models. Overall, three types of experiments were conducted: two at the word level and one at the sentence level.

The novel approach proposed here does not depend on the syntactic structure of the tweets; rather, it extracts semantic and sentimental patterns from a given corpus of tweets.

A. Experiment A: Finding similarities between words

• Objective

The aim was to analyze the semantics and sentiments of tweets at the word level by automatically capturing words with similar semantics and sentiments.

• Method

The unlabeled dataset contained 4,425,003 tweets. A vocabulary of 13,696 tokens was generated and used to create bigram matrix denoted by X, with size = 13,696 x 27,392.

The normalization process was then used. Then tailored bigram matrix was used to discover the most similar words to the query word by comparing between words' vector values in the matrix  $XX^T$ . If words have similar vectors in a matrix, then they tend to have some relatedness. The similarity/relatedness was found by comparing between the vectors (each vector contains 27,392 features) we found the similar words to the query word by comparing which vectors are very closed to the query vectors by using square Euclidean distance function (this matrix was very huge and attempt to open it resulted in "out of memory" problem).

The model arranges all the vocabulary (similar words) in descending order based on similarity with query word (from word that has most highest similar to lower similar word). Then select only the most 10 similar words to make the comparison between models more easy and to make it more clear for reader.

The model was also tested by extracting some sentiment words. The proposed model revealed that words indicative of sentiment tend to have high similarity with words of the same sentiment polarity. The sentiment words do not have different context, the proposed model extract the words that have similarity in polarity or related to query word. If the query word is positive, the model extracts similar/related positive words.

• Results

The results, which were selected and analyzed at random, are shown in tables 2, 3, 4, and 5.

TABLE II. LIST OF THE TEN WORDS MOST SIMILAR TO هلال/HELAL

| Word            | LSA model           | Bigram Model         | Proposed tailored bigram model |
|-----------------|---------------------|----------------------|--------------------------------|
| سعادة/Happiness | راحة/Comfort        | وسعادة/And happiness | راحة/Comfort                   |
|                 | ابديه/Eternal       | سعادة/Happily        | وسعادة/And happiness           |
|                 | جنه/Heaven          | راحة/Comfort         | يسعادة/Happily                 |
|                 | امي/My mother       | سعيدة/Happy          | فرح/Joy                        |
|                 | و فرح/And joy       | راحة/Comfort         | فرحه/Joy                       |
|                 | اكتفاء/Satisfaction | طمائينه/Tranquility  | جنه/Paradise                   |
|                 | اسعدها/Her happy    | و بركه/And blessing  | طمائينه/Tranquility            |
|                 | ~ /Heart symbol     | معالي/Excellency     | حياه/Lifetime                  |
|                 | امنيه/Wish          | عافيه/Good health    | خير/Good                       |
|                 | قلبيها/Her heart    | جنه/Heaven           | رضا/Satisfaction               |

The word sense is set of different meaning of the query word. In the tables 2 and 3 the number refers to how many meanings of word were discovered with respect to usage in text based on the context. In tables 4 and 5 results for sentiment are given. The sentiment words do not have different context, the proposed model extract the words that are similar in polarity or related to query word.

TABLE III. LIST OF TEN WORDS MOST SIMILAR TO هدف/GOAL

| Word            | LSA model           | Bigram Model      | Proposed tailored bigram model |
|-----------------|---------------------|-------------------|--------------------------------|
| هلال/helal      | يقدم/Advent         | يقدم/Advent       | الهلال/Al-helal                |
|                 | شهر/Month           | يا هلال/Ya-helal  | يا هلال/Ya-helal               |
|                 | تهنئه               | ست/Six            | شهر/Month                      |
|                 | /Congratulates      | لشهر/In a month   | بعنه/Expedition                |
|                 | لشهر/For a month    | عتقاء/Redeems     | ست/Six                         |
|                 | سدير/Sadir          | تستقبلون/Greet    | زعيم/Leader                    |
|                 | قدم/Advent          | شهر/Month         | اتحاد/Etihad                   |
|                 | رمضان/Ramadan       | لشهر/For a month  | الغائبين/Absentee              |
|                 | تعذر/Cannot         | بحلول/Advent      | رشيد/Reashed                   |
|                 | مغرد/Twitter's user | لاخر/For the last | سعد/Sad                        |
| المبارك/Blessed |                     |                   |                                |
| Word sense      | 1                   | 1                 | 3                              |

TABLE IV. LIST OF TEN WORDS MOST SIMILAR TO سعادة/HAPPINESS

| Word       | LSA model           | Bigram Model              | Proposed tailored bigram model |
|------------|---------------------|---------------------------|--------------------------------|
| هدف/Goal   | اليوسنه/Bosnia      | قوول/Goool                | قوول/Goool                     |
|            | هدفين/Two goals     | بههدف/By goal or with aim | اهداف/Goals or the aims        |
|            | للمنتخب/To the team | اهداف/Goals or the aims   | الهدف/The goal or the aim      |
|            | المونديال/World Cup | جورول/Goool               | هدفين/Two goals or two aim     |
|            | نيمار/Neymar        | لاعب/Player               | لاعب/Player                    |
|            | التعادل/Equalizer   | تسديده/Win                | بههدف/By goal or with aim      |
|            | الارجنتين/Argentina | تسديده/Score              | جورول/Goool                    |
|            | رونالدو/Ronaldo     | مباراه/Match              | التعادل/Equal                  |
|            | مباراه/Match        | مدرب/Coach                | فوز/Win                        |
|            | كلوزه/Klose         | منتخب/Team                | مباراه/Match                   |
| Word sense | 1                   | 2                         | 2                              |

TABLE V. LIST OF TEN WORDS MOST SIMILAR TO حزن/SADNESS

| Word        | LSA model       | Bigram Model     | Proposed tailored bigram model |
|-------------|-----------------|------------------|--------------------------------|
| حزن/sadness | قلق/Concern     | وحزن/And Sadness | وحزن/And Sadness               |
|             | خدلان/Betrayal  | اشتياق/Longing   | اشتياق/Longing                 |
|             | يكتمل/Complete  | الم/Pain         | حنين/Nostalgia                 |
|             | فقر/Poverty     | فرح/Joy          | الم/Pain                       |
|             | خيبه            | جرح/Wrench       | الم/Longing                    |
|             | /Disappointment | ضيق/Restless     | فرح/Joy                        |
|             | تسكن/Live       | وجع/Wrench       | وجع/Wrench                     |
|             | بلا/Without     | شوق/Longing      | حب/Love                        |
|             | محال/Impossible | الحزن/Sadness    | ضيق/Restless                   |
|             | خيل/Dementia    | النقاء/purity    | جرح/Wrench                     |
| دنيا/World  |                 |                  |                                |

• Discussion

Tables 2 to 5 show the words most similar/related to the given query word by using the proposed tailored bigram, bigram and LSA models. All these methods capture broad

semantic and sentiment relatedness. Based on human evaluation of the experiments results, the tailored bigram model seems to perform better than the LSA model because the proposed model captures more different semantic and sentimental related patterns from a given corpus of tweets.

As can be seen from Table 2, the proposed model shows how the word هلال/helal can have different meanings according to the context; it can mean “crescent,” it can be the name of a Saudi football team, or it can be the name of a person. The word هلال/helal has the meaning “crescent” and is similar to the word شهر/month and the word ست/six, which could be denoted as “number” or “date.” Furthermore, helal, in its meaning as the name of a Saudi football team, is similar to the word زعيم/leader, which is the nickname of the team. The word اتحاد/Etihad is also the name of a Saudi football team. In its meaning as a person’s name, helal is similar to رشيد/Reashed and سعد/Sad, which are also people’s names. The LSA model only gives one semantic context, which is the word helal meaning only “crescent.”

In Table 3, the proposed model shows how the word هدف/goal can have different meanings according to the context; it can mean “score a goal,” and it can mean “aim or target.” Also, the proposed model extracted some informal words, such as فوول/goool and جوول/goool, which are similar to the word هدف/goal. The LSA model only gives one semantic context, which is ‘score a goal.’

The model was also tested by extracting some sentiment words, as shown in tables 4 and 5. The proposed model revealed that words indicative of sentiment tend to have high similarity with words of the same sentiment.

Similarly, in Table 5, the proposed model presents the words فرح/joy and حب/love as being similar to the word حزن/sad—again, all of them have similar words that come after and before then in a sentence. The LSA extracted four words: يكتمل/Complete, تسكن/Live, بلا/Without and دنيا/World were not similar to word sad.

**B. Experiment B: Testing the clustering results**

**• Objective**

The aim was to analyze the semantics of tweets at the word level by automatically capturing patterns of words with similar contextual semantics by using the proposed model (i.e., tailored bigram) which was found to have the highest performance level in the previous experiment.

**• Method**

The train set contained 4,425,003 tweets. A vocabulary of 13,696 tokens was generated and used to create a bigram matrix denoted by X with size = 13696×27392. K-means clustering was then used to categorize the words into k = 100; 200, and 300 clusters. then, made a comparison between k values and arrived at k = 200 as the best result for word clustering to capture patterns of words of similar contextual semantic and sentiment in tweets. Each row was called a vector, and each column was called a dimension (each representing a different semantic feature).

**• Results**

The results are shown in Tables 6, 7 and 8.

TABLE VI. EXTRACTED PATTERNS FOR THE SEMANTIC WORD هلال/HELAL

| Word         | Proposed tailored bigram model |                        |                        |                   |
|--------------|--------------------------------|------------------------|------------------------|-------------------|
|              | Dim 1                          | Dim 2                  | Dim 3                  | Dim 4             |
| هلال / Helal | يقدم /Coming                   | عبدالله /Abdullah      | الغامدي / Al-Ghamdi    | تشكيله /Formation |
|              | بحلول /Advent                  | سعد /Sad               | الحربي / Al-Harbie     | لاعبة /Players    |
|              | بشهر /At month                 | عبدالرحمن /Abdulrahman | القحطاني / Al-Qahtani  | جماهير /Masses    |
|              | بمناسيه /Occasion              | فهد /Fahd              | الشمري /AlShammari     | مباراه /Match     |
|              | بالعيد /Eid                    | خالد /Khalid           | الدوسري / Al-Dosari    | فوز /Win          |
|              | يقرب /Near                     | بندر /Bender           | الانزي /AL-Anzi        | رمي /Goal         |
|              | برمضان /Ramadan                | ناصر /Nasser           | بن /Son of             | يشجع /Encourages  |
|              | قدم /Coming                    | فيصل /Faisal           | الشهراني / Al-Shahrani | مدرب /Coach       |

TABLE VII. EXTRACTED PATTERNS FOR SEMANTIC WORD/هدف GOAL

| Word      | Proposed tailored bigram model |                  |                 |                      |
|-----------|--------------------------------|------------------|-----------------|----------------------|
|           | Dim 1                          | Dim 2            | Dim 3           | Dim 4                |
| هدف /Goal | تشكيله /Formation              | نيمار /Neymar    | قرار /Decision  | بلجيكا /Belgium      |
|           | لاعبة /Players                 | سواريز /Suarez   | طلب /Request    | هولندا /Netherlands  |
|           | جماهير /Masses                 | ميسي /Messi      | حساب /Account   | الارجنتين /Argentina |
|           | مباراه /Match                  | سانشيز /Sanchez  | موضوع /Topic    | المانيا /Germany     |
|           | فوز /Win                       | كوستا /Costa     | اخبار /News     | تشيلي /Chile         |
|           | رمي /Goal                      | بنزيما /Benzema  | موقع /Position  | فرنسا /France        |
|           | يشجع /Encourages               | رونني /Ronnie    | اسم /Name       | غانا /Ghana          |
|           | مدرب /Coach                    | رونالدو /Ronaldo | بيان /Statement | كولومبيا /Columbia   |

TABLE VIII. EXTRACTED PATTERNS FOR THE SEMANTIC WORD خطير/DANGEROUS

| Word            | Proposed tailored bigram model |                    |                  |                  |
|-----------------|--------------------------------|--------------------|------------------|------------------|
|                 | Dim 1                          | Dim 2              | Dim 3            | Dim 4            |
| خطير /Dangerous | رائع /Wonderful                | هلالي /Hilali      | محظوظ /Lucky     | موجع /Painful    |
|                 | جيد /Good                      | برازيلي /Brazilian | محترم /Respected | مؤلم /Painful    |
|                 | مميز /Special                  | نصراوي /Nasraoui   | مبسوط /Happy     | المؤلم /Painful  |
|                 | متنع /Fun                      | السعودي /Saudi     | اناني /Selfish   | يوجع /Pain       |
|                 | سيء /Bad                       | اهلاوي /Ahlawy     | غلطان /mistaken  | يؤلمني /Pains me |
|                 | جميل /Beautiful                | عربي /Arabian      | كذاب /Liar       | يوجعني /Pains me |
|                 | ممتاز /excellent               | مريدي /Madrída     | مظلوم /oppressed | يوجعك /Pain      |
|                 | كبير /Big                      | مصري /Egyptian     | مجنون /crazy     | قاسي /tough      |

Tables 6 to 8 give the results for the four most common dimensions for the given query words using the proposed model after reduction of the dimensions to 200 features. All these query words captured broad contextual semantic

similarities. K-means clustering was used to determine which words belonged to each cluster.

Table 6 presents the top four semantic features (or dimensions) for the word هلال/helal. The first dimension indicates the results obtained from mining the meaning “crescent.” The second dimension is related to the word’s meaning as a person’s name, and the third dimension indicates the word’s meaning as a family name. The fourth dimension is connected with the meaning of helal as the name of a Saudi football team.

Table 7 presents the key semantic features of the word goal. The first dimension indicates the results obtained from mining the meaning “score a goal.” In the second dimension, the word indicates the names of football players, which are also in the sport domain. In the third dimension, the word gives the meaning of the “aim or target”. The fourth dimension connects the word with its meaning in relation to the names of countries. The words and phrases goal, player, team name, and country name were all found before the word سجل /score, explaining why these words appeared in these particular dimensions.

In table 8 the four most common semantic features for the word dangerous are presented. The first dimension refers to the word’s meaning as “wonderful” (i.e., positive). The second dimension is not related to the word dangerous. The third dimension indicates the first three results gave the word dangerous meaning wonderful (i.e., positive), and the last five results gave the word dangerous meaning bad (i.e., negative). The fourth dimension refers to the word’s meaning as dangerous (i.e., negative).

The proposed model categorizes words together that have similar semantic features, automatically capturing the contextual patterns in tweets. If a word has multiple contextual meanings, the model uncovers these meanings, adding each word to the relevant clusters.

C. Experiment C: Finding similarities at sentence level

• Objective

The aim is to analyze the semantics and sentiments of tweets at the sentence level by automatically capturing sentences with similar semantics and sentiments. The words هلال/helal, </goal>/هدف, and sentiment word خطير/Dangerous were tested in the results for the proposed model at the sentence level.

• Method

First, all the sentences in the database that contained the query words were extracted. The vector average of each sentence was then calculated using the bigram matrix X after reducing the dimensions as input in equation (2), thus giving a new matrix, V si×200. The dimensions were reduced because the bigram matrix used in calculating matrix V caused a memory problem with the computer. The reduction eliminated the problem

The five sentences that were most similar to each query sentence were tested by comparing the sentence vectors using matrix V.

• Results

The results are shown in tables 9, 10, 11, 12, 13, 14 and 15

TABLE IX. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD هدف/GOAL, WHERE THE WORD هدف MEANS “GOAL”

| Original Sentence  | Most similar sentences   |
|--|--|
| عندما يكون لديك هدف لا تساط تركيزك على المستقل وتخسر ، سلط الان الضوء على وهي حتما الان ستصلك الي ما تريد  | لا يوجد شخص ولد كبيرا .. ولا يوجد مشهور لم يبدأ صغيرا .. ولا يوجد هدف تحقق الا وكان حلمًا وليدالبدايه من الصفر ليست عيبا<br>“No person is born an adult. No famous person does not start small. There is no goal archived, until there was a dream that born from scratch, this is not a defect”   |
| “When you have a goal, do not fix your focus on the future and lose the ‘now’; highlighted the ‘now’ and, inevitably, you will receive what you want.” | : لا تعيش بلا امل ولا هدف , ولا تقول ان الزمن @alfofonumber دايم بخيالجمال الاشياء تلقاها صدف .. والفرح دايم وري الصبر الجميل !!<br>“@alfofonumber: do not live without hope nor objective, and do not say that time is always miserly; make things beautiful, received coincidence, and find joy always behind beautiful patience!!”  |
|  | حدد هدفك واستعن بالله ببساطه كم تريد ان تختم كتاب الله في رمضان #الاستعداد لرمضان #رمضان<br>“Simply how many times you want to seal the Book of Allah in Ramadan? Determine your goal and seek the help of God #prepare for Ramadan #Ramadan.”   |
|  | @ zlfay qnumbermeNUMBER @ aNUMBERm@ abnazulfi@ والشكر ومن تقاغل وهدفتا جميعا المصلحه العامه حفظ الله الوطن والقائمين عليه متمنين ان تحل هذه المشكله<br>“@QNUMBERmeNUMBER @aNUMBERm @abnazulfi @zlfay and thanks to you also for how interact, and our aim of all is the public interest. God Save the homeland and those who support it and wish that interaction would solve this problem.” |
|  | صوره مؤثره شاب سوري يطعم طفله من تحت باب منزلها بهدف طمأنتها والتخفيف عنها<br>An image of an inuential young Syrian feeding a child from under the door of her home in order to reassure her and to comfort her  |

TABLE X. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD هدف/GOAL, WHERE THE WORD هدف MEANS “AIM”

| Original Sentence  | Most similar sentences   |
|--|--|
| يسألوني من عقب الهلال تحب؟! # الهلال يايعدي مايعد ه احد                                | “Ask me after the al-helal love!? # al-helal my love, no one came after it.”   |
| تتواجد عدد من الجماهير العمانيه بمقر نادي # الهلال جاءت لتوازر الفريق الازرق امام السد | “Number of Omani fans present at the Club #al-helal came to co-operate the blue team in front of al-Sad. URL.”   |
| مبروك بطوله الليومزين URL يا هلال  | سنوات فتح الغرافه الملعب NUMBER قبل @bluegoldNUMBER بالكامل لجمهور الهلال .. لانها الخلاق الكبير ولا URL<br>“@ bluegoldNUMBER: NUMBER years before Gharafa opens the entire stadium to the al-Helal audience. because this is manners of the greater no consolation to Abu five .. !! URL.”  |
| “Congratulations Championship Limousine O Helal URL”                                   | لو الهلال حقق اسيا وش بتسويفي حال اقام الاتحاد الاسيوي بطوله اسيا # بمشاركة الهلال وحده فقط ، اقول ربما يحقق البطولة .. اقول ربما “If al-Helal achieved the champion of Asia, what you will do in case if the AFC established Asian Cup with the participation of al-Helal alone only, say perhaps he will achieve tournament .... I side probably.” |
|  | لما تكون نادي كبير والرغم الاول جماهيرا ، فانت تحتاج الى مركز اعلامي @ a bin mosaad قوي جدا !! وما يحدث في نادي # الهلال هو العكس تماما<br>“When you become a big club and the first number to be a mass, you need a media center that is very strong!! What happens in Club #al-Helal is quite the opposite! @abinmosaad”                           |

TABLE XI. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN WORD هلال/HELAL, WHERE THE WORD هلال MEANS "THE NAME OF A SAUDI FOOTBALL TEAM"

| Original Sentence  | Most similar sentences   |
|--|--|
| <p>ميروك بطوله الليموزين يا هلال URL</p> <p>"Congratulations Championship Limousine O Helal URL"</p> | <p>يسألوني من عقب الهلال تحب؟! # الهلال يا بعدى مابعد ه احد</p> <p>"Ask me after the al-helal love! # al-helal my love, no one came after it."</p>   |
|  | <p>تتواجد عدد من الجماهير العمانيه بمقر نادي # الهلال جاءت URL.</p> <p>"Number of Omani fans present at the Club #al-helal came to co-operate the blue team in front of al-Sad. URL."</p>  |
|  | <p>سنوات فتح الغرافه NUMBER قبل NUMBER @bluegoldNUMBER</p> <p>الملاعب بالكامل لجمهور الهلال .. لانها اخلاق الكبير ولا URL</p> <p>"@bluegoldNUMBER: NUMBER years before Gharafa opens the entire stadium to the al-Helal audience, because this is manners of the greater no consolation to Abu five .. !! URL."</p>  |
|  | <p>لو الهلال حقق اسيا وش بتسوييني حال اقام الاتحاد الاسيوي # بطوله اسيا بمشاركة الهلال وحده فقط ، اقول ربما يحقق البطولة .. اقول ربما</p> <p>"If al-Helal achieved the champion of Asia, what you will do in case if the AFC established Asian Cup with the participation of al-Helal alone only, say perhaps he will achieve tournament ... I side probably."</p> |
|  | <p>لما تكون نادي كبير والرقم الاول جماهيرا ، فانت تحتاج الى مركز اعلامي قوي جدا !! وما يحدث في نادي # الهلال هو @ a bin mosaad ! العكس تماما</p> <p>"When you become a big club and the first number to be a mass, you need a media center that is very strong!! What happens in Club #al-Helal is quite the opposite! @abinmosaad"</p>                            |

TABLE XII. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD هلال/HELAL, WHERE THE WORD هلال MEANS "THE NAME OF A PERSON"

| Original Sentence   | Most similar sentences   |
|---|--|
| <p>وابع من العيش ما تسر به ان عتل الناس فيه او عنروابو هلال العسكري</p> <p>"I want from the living what pleased, if Conquer people with it Abu Helal al-Askarry."</p> | <p>الشهادة لله ان نفتخر بالشباب الي مثل اخي ابو @alialnimi هلال سباق لفعل الخير جزاك الله خيرا . ونحن في الخدمة @TheNaim testimony to God that we are proud to youth like my brother Abu Helal, he racing to do good things, God reward you. We are in the service any time."</p>  |
|   | <p>بين ما يقدرون لانهم بيتابعون mbc الهلاليين بيون يقاطعون السعودية الاولى</p> <p>"Alhlaliyn wants to boycott MBC but they cannot because they follow first Saudi channel."</p>  |
|   | <p>تطرح شركة صله تناكر مباراه # الهلال والعروبه اليوم بالنادي مساء NUMBER عصرا حتى NUMBER من الساعه ويتواصل بيع التناكر غدا بالملاعب باستاد الملك فهد</p> <p>"The Sela company raises the number of tickets for the #Al-Helal and Arabism match today from NUMBER pm until NUMBER evening and will continue selling tickets tomorrow at King Fahd Stadium"</p> |
|   | <p>متى يلعب الهلال واشوف التايم لاين كله ازرق والجمهور يملى مدرجات الهلال ونشوف اللاعبين متى يجي اليوم اللي اكون مبسوطة كثير</p> <p>"When al-Helal playing and see the whole twitter favorite is blue and the al-Helal's fans full the stadium, and see players</p> <p>When the day comes when I can be joyous.</p>  |
|   | <p>سبق اني ارسلت تغريده تخصص ديناه الهلال مصدرها صحيفه سيورت الاكترونيه وبعد ان وضع ان ما كتب غير صحيح اترفع باخلاقي كمسلم نصراوي واعتذر</p> <p>"I already sent a tweet about al-Helal religion source electronic newspaper Sport and after it explained that what has been written is not true, as Muslim manners and Nasraoui, I apologized."</p>            |

TABLE XIII. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD هلال/HELAL, WHERE THE WORD هلال MEANS "CRESCENT"

| Original Sentence  | Most similar sentences  |
|--|---|
| <p>عاجل المشروع الاسلامي لرصد الاهله : تمت رويه هلال شهر شوال قبل قليل نهارا في السعودية باستخدام تقنيه التصوير الفلكي #هلال شوال رمضان</p> <p>"Urgent Islamic Crescents Observation Project: This sighting of the new moon of Shawwal shortly before daybreak in Saudi Arabia uses the technique astrophotography #crescentofShawwal #Ramadan."</p> | <p>كل العالم ينتظر هلال الفطر ( شوال ) الا اهل عزه فانهم ينتظرون هلال النصر ويقولون متى هو قلى عسى ان يكون قريب</p> <p>"All the world is waiting for the crescent al-Fitr (Shawwal), but the people of Gaza, they are waiting for the crescent's win and say when it Tell It may be that close."</p>                              |
|  | <p>@balsayegh @ arabicobama الجهات الحكوميه بالسعوديه تتبع تقويم ام القرى وليس رويه الهلال</p> <p>"@Balsayegh @arabicobama government agencies in Saudi Arabia follow the calendar or Imm Alqri; they do not see the crescent."</p>   |
|  | <p>الاتحاد سجل هدف ومنتر الهلال وتلعب الكوره في نصف عرفنو كيف URL الملعب وبعد دقائق يلغى الحكم هدف الاتحاد</p> <p>NUMBER البطولة = "Al-Etheaad scored and al-Helal's center and play a ball in half pitch and after minutes the ruling canceled Al-Etheaad goal URL know now how the championship = NUMBER."</p>                  |
|  | <p>بامكانك تحديد حجم عقليه المشجع الهلالي من خلال حجم طاقيته - طاقيه صغيره - يقولك " السابجه تقرب " - طاقيه كبيره - يقولك " الهلال ملكي "</p> <p>"You can determine the size of the mentality encouraging Hilali through the cap size—small cap tells you 'seven o'clock approaching'; large cap tells you 'Royal al-Helal.'"</p> |
|  | <p>انا اموت ولا اغير الهلال @eeNUMBERqwe</p> <p>"@eeNUMBERqwe If I die, it does not change my al-Helal."</p>  |

TABLE XIV. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD DANGEROUS, WHERE THE WORD DANGEROUS MEANS "DANGEROUS"

| Original Sentence   | Most similar sentences  |
|---|---|
| <p>طرق لشرب NUMBER الشاي تجعله خطير جدا URL على صحتك</p> <p>NUMBER ways to drink tea to make it very dangerous to your health URL</p> | <p>جالس اتابع مباراه العين والنصر اعاده اجانب العين ماقيهم الا جيان الخطير اما كيمبو والسولفاكي والكوري مستواهم عادي الحمد لله</p> <p>I am following the rematch between Al-Aean and al-Nasser, at al-Aean, the foreign player Jian is dangerous but Kimpo and Alsohgaki and Korea are normal</p>   |
|   | <p>دنبلي خطير جداوسريع Denbla very dangerous and fast</p>   |
|   | <p>فتوى ان من اجاز الاغاني مجاهرالاتجوز امامته رايب انها زله خطيره</p> <p>The fatwa of authorized songs boldness may not be Imamth my opinion it slip serious</p>   |
|   | <p>@ soumahran ايه يابنتي جيتي المستندات الخطيره دي منين</p> <p>hi from where you got this Serious Documents</p>  |
|   | <p>عباره خطيره تؤدى الي الشرك وهي ان تقول [ بكره يجلها الف حلال ] والصح [ يجلها الله سبحانه ] فقول الف حلال تعني ان " هناك الف رب " لاله الا الله</p> <p>Dangerous word lead to polytheism which is to say [tomorrow will solve it thousand solver] the correct say [solved by God] To say a thousand solver means that there are a thousand Lord "to God but God."</p> |

TABLE XV. THE FIVE MOST SIMILAR SENTENCES THAT CONTAIN THE WORD DANGEROUS, WHERE THE WORD DANGEROUS MEANS “WONDERFUL”

| Original Sentence  | Most similar sentences   |
|--|--|
| @ fahdalruqi خطر الله ابو عمر و شاعر بعد<br>@ fahdalruqi God you are wonderful Abu Omar, and you also poet | @ tntn NUMBER تابعيها مره خطيره<br>@ tntn NUMBER follow her she is wow   |
|  | اصابه نيمار الخطيره في مباراه البرازيل وكولومبيا .. كسر في فقره الظهر خطيره جدا URL<br>Neymar had serious injury in the match between Brazil and Colombia .. a broken vertebra back very serious URL |
|  | انه امر خطير جدا استقبال النكت باستمرار خطير على تكوين ابناؤنا وبناتنا<br>It's too dangerous receiver jokes constantly risk of the formation of our sons and daughters                               |
|  | بانزيما خطير والله خطير قول<br>Benzema goal is grave and serious   |
|  | @ noda NUMBER الصورة خطيره ☹<br>@ Noda NUMBER Image is serious ☹   |

### • Discussion

Tables 9 to 15 give the five most similar sentences to the given query sentences using the word representations generated by the proposed model. All these vectors capture broad semantic and indirect sentiment similarities.

Table 9 gives the five sentences containing the word goal that are most similar to the query sentence where the word goal means “scored a goal.”

Table 10 gives the five sentences containing the word goal that are most similar to the query sentence, where the word goal means “aim” or “target.” The proposed model extracts similar semantic contextual sentences that contain the word goal where it means “aim” or “target.” The results in table 8 and 9 show the two different semantic contexts for the word goal at the sentence level.

Table 11 presents the five sentences containing helal that are most similar to the query sentence, where the word helal means “the name of a Saudi football team.” All the similar sentences contain helal where it denotes the name of a Saudi football team.

Table 12 gives the sentences where helal means “the name of a person.” The proposed model only extracts the one sentence where helal means “a person’s name” in the sentence context. The other similar sentences refer to helal as the name of a Saudi football team and are not similar to the query sentence.

Table 13 presents the sentences where helal means “crescent.” Here, the model has only extracted two sentences where the word helal means “crescent;” the other three are examples where the word helal refers to the Saudi football team and are not similar to the query sentence. The results in table 11, 12, and 13 show the three different semantic contexts for the word helal at the sentence level.

Table 14 presents the sentences where Dangerous, where Dangerous means “Dangerous”. The model has extracted the five sentences containing the word Dangerous that are similar to the query sentence, where Dangerous means “Dangerous”. Table 15 presents the sentences where Dangerous means “Wonderful”. Here, the model has only extracted gives the three sentences containing the word Dangerous that are similar

to the query sentence, where Dangerous means “Wonderful”. The results in table 14 and 15 show the two different semantic contexts for the word Dangerous at the sentence level.

The proposed model has been used to analyze the semantics and sentiments of tweets at the sentence level, automatically capturing the patterns of sentences with similar contextual semantics and sentiments in tweets. According to the model results, the method needs to be developed further in order for more accurate results to be obtained.

### V. CONCLUSION

The proposed tailored bigram model used unsupervised clustering at word and sentence level to allow semantic and sentiment categorization to take place. In the experiments, words and sentences in tweets with similar semantics and sentiments were automatically captured and grouped. The proposed model was then compared with the classic bigram and LSA models. The proposed approach was not concerned with the syntactic structure of tweets, but with the extraction of patterns in semantics and sentiments from a particular tweet corpus.

With this methodology, a huge corpus was used, no annotation processing was utilized for labels, the word order within the tweets was considered, and no filtering process was used. The filtering was used only to “clean” the text, thus reducing the corpus size and the noise in the text. These steps were taken to ensure that the contexts of the tweets remained unchanged. Semantic dictionaries or lexicons were not used due to their limited coverage for informal Arabic. Based on the work, we conclude that although difficult to handle, big data can help in checking almost every type of possibility of similarity/ relatedness among words. Although due to availability of limited computational resources, we used some threshold to reduce the data, but were still were able to get good results. The manual evaluations of the results need to be automated for which Arabic semantic resources should be developed.

### ACKNOWLEDGMENT

This work was supported by King Abdulaziz City for Science and Technology (KACST) Project number AT-200-34

### REFERENCES

- [1] L. Albraheem and H. S. Al-Khalifa, “Exploring the problems of sentiment analysis in informal Arabic,” in *Proceedings of the 14th international conference on information integration and web-based applications and services*, 2012, pp. 415–418.
- [2] A. Shoukry and A. Rafea, “Sentence-level Arabic sentiment analysis,” in 2012 international conference on collaboration technologies and systems (CTS), 21-25 May 2012, 2012, pp. 546–550.
- [3] B. Liu, “Sentiment analysis and opinion mining,” *Synth. Lect. Hum. Lang. Technol.*, vol. 5, no. 1, pp. 1–167, 2012.
- [4] R. T. Khasawneh, H. A. Wahsheh, M. N. Al-Kabi, and I. M. Alsmadi, “Sentiment analysis of Arabic social media content: A comparative study,” in 2013 international conference on Information Science and Technology (ICIST), 9-12 Dec. 2013, 2013, pp. 101–106.
- [5] M. Rushdi-Saleh, M. T. Martín-Valdivia, L. A. Ureña-López, and J. M. Perea-Ortega, “OCA: Opinion corpus for Arabic,” *J. Am. Soc. Inf. Sci. Technol.*, vol. 62, no. 10, pp. pp. 2045–2054, 2011.
- [6] H. Saif, Y. He, M. Fernandez, and H. Alani, “Semantic patterns for sentiment analysis of Twitter,” *Semant. Web-ISWC 2014*, vol. (Vol. 8797, pp. pp. 324–340), 2014.

- [7] C. C. Aggarwal and C. Zhai, Mining text data. Berlin & Heidelberg, Germany: Springer, 2012.
- [8] C. Zhai, "Statistical language models for information retrieval," *Synth. Lect. Hum. Lang. Technol.*, vol. 1, no. 1, pp. 1–141, 2008.
- [9] J. Lin and C. Dyer, "Data-intensive text processing with MapReduce," *Synth. Lect. Hum. Lang. Technol.*, vol. 3, no. 1, pp. 1–177, 2010.
- [10] M. Moussa, M. W. Fakh, and K. Darwish, "Statistical denormalization for Arabic text," in *Empirical Methods in Natural Language Processing*, 2012, vol. 228, pp. 228–232.
- [11] W. Naptali, M. Tsuchiya, and S. Nakagawa, "Word co-occurrence matrix and context dependent class in LSA based language model for speech recognition," *Int. J. Comput.*, vol. 3, no. 1, pp. 1–11, 2009.
- [12] D. Laniado and P. Mika, "Making sense of Twitter," in *The Semantic Web–ISWC 2010*, vol. 6496, P. F. Patel-Schneider, Y. Pan, P. Hitzler, P. Mika, L. Zhang, J. Z. Pan, I. Horrocks, and B. Glimm, Eds. Berlin & Heidelberg, Germany: Springer, 2010, pp. 470–485.
- [13] R. M. Duwairi, R. Marji, N. Sha'ban, and S. Rushaidat, "Sentiment Analysis in Arabic tweets," in *5th international conference on information and communication systems (ICICS)*, 1-3 April 2014, 2014, pp. 1–6.
- [14] A. E.-D. A. Hamouda and F. E. El-Taher, "Sentiment analyzer for Arabic Comments System," *Int. J. Adv. Comput. Sci. Appl.*, vol. 4, no. 3, pp. 99–103, 2013.
- [15] M. M. Itani, L. Hamandi, R. N. Zantout, and I. Elkabani, "Classifying sentiment in Arabic social networks: Naive Search versus Naive Bayes," in *2012 2nd international conference on advances in computational tools for engineering applications (ACTEA)*, 12-15 Dec. 2012, 2012, pp. 192–197.
- [16] S. R. El-Beltagy and A. Ali, "Open issues in the sentiment analysis of arabic social media: A case study," in *2013 9th international conference on innovations in Information Technology (IIT)*, 17-19 March 2013, 2013, pp. 215–220.
- [17] N. A. Abdulla, N. A. Ahmed, M. A. Shehab, and M. Al-Ayyoub, "Arabic sentiment analysis: Lexicon-based and corpus-based," in *2013 IEEE Jordan conference on applied electrical engineering and computing technologies (AEECT)*, 3-5 Dec. 2013, 2013, pp. 1–6.
- [18] M. Taboada, J. Brooke, M. Tofiloski, K. Voll, and M. Stede, "Lexicon-based methods for sentiment analysis," *Comput. Linguist.*, vol. 37, no. 2, pp. 267–307, 2011.
- [19] H. Saif, Y. He, and H. Alani, "Semantic sentiment analysis of Twitter," in *The Semantic Web–ISWC 2012*, vol. 7649, P. Cudré-Mauroux, J. Heflin, E. Sirin, T. Tudorache, J. Euzenat, M. Hauswirth, J. X. Parreira, J. Hendler, G. Schreiber, A. Bernstein, and E. Blomqvist, Eds. Berlin & Heidelberg, Germany: Springer, 2012, pp. 508–524.
- [20] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proceedings of the 49th annual meeting of the Association for Computational Linguistics: human language technologies*, 2011, vol. 1, pp. 142–150.
- [21] H. Froud, A. Lachkar, and S. A. Ouatic, "Arabic text summarization based on latent semantic analysis to enhance Arabic documents clustering," *Int. J. Data Min. Knowl. Manag. Process*, vol. 3, no. 1, pp. 79–95, 2013.
- [22] H. M. Alghamdi, A. Selamat, and N. S. A. Karim, "Arabic web pages clustering and annotation using semantic class features," *J. King Saud Univ. Inf. Sci.*, vol. 26, no. 4, pp. 388–397, 2014.
- [23] R. Sarikaya, M. Afify, Y. Deng, H. Erdogan, and Y. Gao, "Joint morphological-lexical language modeling for processing morphologically rich languages with application to dialectal Arabic," *IEEE Trans. Audio. Speech. Lang. Processing*, vol. 16, no. 7, pp. 1330–1339, 2008.
- [24] A. E.-D. Mousa, R. Schluter, and H. Ney, "Investigations on the use of morpheme level features in language models for Arabic LVCSR," in *2012 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 25-30 March 2012, 2012, pp. 5021–5024.
- [25] H. Froud, A. Lachkar, and S. A. Ouatic, "Stemming for Arabic words' similarity measures based on Latent Semantic Analysis model," in *2012 international conference on multimedia computing and systems (ICMCS)*, 10-12 May 2012, 2012, pp. 779–783.
- [26] H. Froud, A. Lachkar, and S. A. Ouatic, "Stemming versus Light Stemming for measuring the similarity between Arabic Words with Latent Semantic Analysis model," in *2012 colloquium in Information Science and Technology (CIST)*, 22-24 Oct. 2012, 2012, pp. 69–73.
- [27] N. M. Zaki, K. A. Alawar, A. A. Al Dhaheri, and S. Harous, "Protein-protein Interaction Prediction using Arabic semantic analysis," in *2013 9th international conference on innovations in Information Technology (IIT)*, 17-19 March 2013, 2013, pp. 243–247.