

Sentiment Analysis of Social Media Micro Blogs using Power Links and Genetic Algorithms

Mahmoud Rokaya

College of Computers and Information Technology, Taif University
Taif, Saudi Arabia
Faculty of Science, Tanta University
Tanta, Egypt

ABSTRACT

In the current years, social media became one of the most important sources of data for different data analytic purposes. One of the most important issues is how to map different trends in social media and define the relation between different groups based on their sentiment or interests. In this paper, a two-phase approach is used for clustering a set of blogs. At the first phase, the approach builds a lexicon that provides the polarity of each word. In the second phase, the approach clusters the blogs bases on the polarity features and the Power Link features. The output of the second phase is used as the input of the first phase to get an improved lexicon. This process will continue in a loop between phase one and phase 2 till a stable set of clusters is gotten. The approach aims to develop a non-supervised cluster agent that can correctly cluster micro blogs and define different interests of different groups of people. The results of the approach are expressed in terms of precision, recall and F-measure.

General Terms

Artificial Intelligence. Natural Language Processing. And Text Analysis.

Keywords

Sentiment Analysis, Evolution Calculation, Genetic Algorithms, Power Links, Social Media, Micro Blogs.

1. INTRODUCTION

Classifications and clustering are two basic tasks in machine learning and data science [1]. Classifications are used when a set of labels are known, and it is needed to define the suitable label for a specific item [2]. Supervised learning is related to classification problem. In the supervised learning, a set of items with known labels and this set is used to train the machine hoping that it can give the write label for items that their labels are unknown. [3]. Clustering means to group similar items in groups, there this no known labels in advance [4]. Clustering is related to unsupervised learning algorithms. In the clustering problem, there is a set of items and it is needed to group each similar set of items in one group. In classifying and clustering, the concept of features and extraction of the features are essential tasks [5]. In classifying, the training aims to train the machine to combine the features in some way to decide to which label the item best probably could be mapped. In clustering, the machine depends on the features and the distance function. Distance function computes the similarity between items using the distance function. Similar items will have shorter distance [6].

In social media sentiment analysis, using clustering and classification is popular [7]. The classifications approaches can be Lexicon-based and corpus-based methods [8] The

approach of Lexicon-based methods depends on predefined static lexicons to determine the polarity of individual words in a corpus, then using a polarity function to determine the polarity of each blog. There are many works that adopted this approach, for example [9]. The corpus-based approach depends on the semantic relation of the words. This make the semantic approach is language dependent [7]. There are many works that adopted this approach, for example [10, 11]. The second approach is language dependent and the first approach is independent from language. Also, dependency of the lexical approach on a fixed lexicon seems to be unsuitable with the highly speed chaining of the content in social media. Also, depending on the language properties seems to be not much effective, considering that the social media users rarely respect the grammar or spelling of the language , also, many times, it is found that the users uses mixed languages or even used the alphabet of another language to express their own words. However, in terms of accuracy the semantic approach gave better results. Cluster methods used in social media are metadata-based approaches, lexical semantics approaches or hybrid. The lexical semantic approaches depend on the meta data of the blogs not on the words on the words in the blog. From the other hand, the contextual semantics approaches depend on the context properties of the blogs [12]. For example [13, 14] used the external semantic properties to cluster hashtags. [15, 16] used the contextual sematic properties. [12] used a hybrid approach.

The mentioned methods either use the supervised learning, which become not preferred these millions of tons of data that appear every day, or they are related to the language features or to sematic features of the content or the meta-sematic and this means that it is needed to develop a different method for each group of languages. Depending on quantitative data a clustering method will solve the two problems. For sure adding the languages features or the meta-structure of the language will improve the performance. In this paper, a clustering method for micro blogs that is not directly dependent on language properties will be presented.

Genetic Algorithm has its origin as a method to solve the optimization problem regardless the nature of the optimization function or the constrains [17]. Genetic algorithms adopt the natural evolution [18]. [7,8] used the genetic algorithms to optimize lexicons for sentiment analysis. [8] explained the uses of genetic algorithms uses in general and for specific purpose of sentiment analysis and building lexicons.

A lexicon is the vocabulary of a person, language, or branch of knowledge. Lexicons in most cases are represented by two columns, the first column, presents the words, the second columns present the function of the lexicon. For sentiment analysis, the second column presents the polarity value of the

word. Some works tried to build lexicons based on solving an optimization problem such as [7, 8].

Power Links concept was suggested by Rokaya and Atlam [19]. Power links has many applications in text extraction and summarization [20, 21]. Power links was used in context spelling [22, 23] and information retrieval [24].

This paper aims to develop a method for clustering micro blogs. The method is dependent from the language's properties, so it can be applied for any language. The method adopts a non-supervised approach. The method can be divided into two main phases, the lexicon building phase and the clustering phase. The output of the first phase is the output for the second phase, so, there is a loop that can be repeated until reaching accepted results. In the second phase, the clustering process goes in two steps, the first step is a clustering based on Power Links. Each cluster resulting from the power link clustering step is further classified into two groups based on the polarity of each blog the cluster. The polarity is calculated based on the polarity of the words. Polarity of words is gotten from the generated lexicon in the first phase.

The remaining sections of the paper are as follows. Section 2 presents the general architecture of the proposed method. It also gives the basic definitions of Power Link (adopted in this work) and the genetic algorithm general architecture. Section 2 gives the details of building the lexicon based on the genetic algorithm and the details of the clustering algorithms. Section 3 presents the experiments that are used to test the method. Finally, section 4 presents the results and the future work.

2. THE METHOD

Fig1 shows the architecture of the clustering algorithm. The first phase used to produce the lexicon. The second phase uses the lexicon beside the Power Link features to produce the clusters of the blogs. Again, these clusters will be used as the new base to rebuild the lexicon and again the lexicon will be used beside the Power Links to produce the new clusters. This process will be repeated until the clusters become stable. The initial lexicon is any available one. The method depends on developing the ability of the method through improving the lexicon through training loops as well as adding new words of power links that contribute in giving the correct class of a blog. The method depends on considering the current class are the base to classify a new set of blogs. Adding the new set will modify the values of the polarity of each word in the current lexicon to improve the lexicon ability to classify new blogs with a higher accuracy. This process will continue and by the time the chance of improving the method performance will never stop. The experiments will show how this process can proceed.

Let w_i be a word in a blog b_j from a set of blogs S_k . Suppose a lexicon l_m is used and the polarity score of w_i in l_m is given by $pwl(w_i)$ then the polarity score of b_j is given by:

$$PSB(b_j, S_k, l_m) = \sum_{w_i \in b_j} pwl(w_i)$$

To decide the polarity class of a blog b_j , the function $PC(b_j, l_m)$ is used

$$PC(b_j, l_m) = \begin{cases} R1_+ & \text{if } PSB(b_j, S_k, l_m) > 0 \\ R1_- & \text{if } PSB(b_j, S_k, l_m) \leq 0 \end{cases}$$

The accuracy of l_m is given its ability to map blogs to the correct class. Let $A(l_m, S_k)$ be the accuracy function, then

$$A(l_m, S_k) = \frac{\text{number of blogs in } S_k \text{ that } l_m \text{ predicted their class correctly}}{\text{number of blogs in } S_k}$$

Now, one can write the problem of finding the best l_m as an optimization problem.

$$\text{For } S_k, l_{opt} = \underset{l_m}{\text{argmax}} A(l_m, S_k)$$

According to [8], it is better to use a penalty function as a fitness function rather than using the accuracy function itself. The penalty function can be stated as:

$$P(S_k, b_j, l_m) = \begin{cases} \frac{|PSB(b_j, S_k, l_m)|}{\theta} & \text{if } PSB(b_j, S_k, l_m) \text{ failed to classify } b_j \text{ correctly} \\ -\frac{|PSB(b_j, S_k, l_m)|}{\theta} & \text{if } PSB(b_j, S_k, l_m) \text{ classified } b_j \text{ correctly} \end{cases}$$

Where θ is the penalty factor. The lexicon will get a positive penalty if it failed to classify a blog to its correct class and it will get a negative penalty if the lexicon succeeded to classify a given blog to its correct class. The fitness function for a lexicon l can be written as the sum of penalty values for each blog classifies using l . Fig. 2 illustrates the steps of how to use the penalty $P(S_k, b_j, l)$ function as a fitness function $F(S_k, l) = -\sum_{b_j \in S_j} P(S_k, b_j, l)$

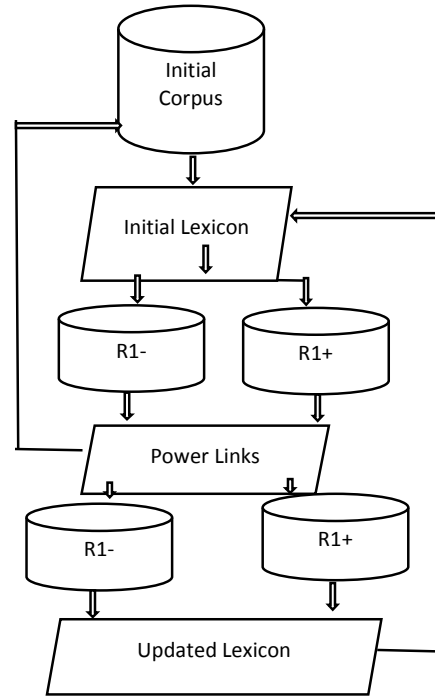


Fig.1 System Overview

The genetic algorithm works through five basic steps, initialization, selection, crossover, mutation and replacement. The initialization step means to initiate the values of all chromosomes randomly. In this case the chromosomes are available lexicons in a given range. Two values, maxpol and minpol are chosen then the polarity of each word is set randomly to be any integer between maxpol and minpol. The selection step is done through the random selection based on roulette wheel strategy. The crossover is achieved by replacing the values in two chromosomes in predefined places with a specific probability. Also the mutation is achieved by

selecting some chromosomes in random places in the same chromosome and for a specific probability. The fitness function determines which set of chromosomes will be

selected in the next evolution. Fig.3 shows the detail of the algorithm

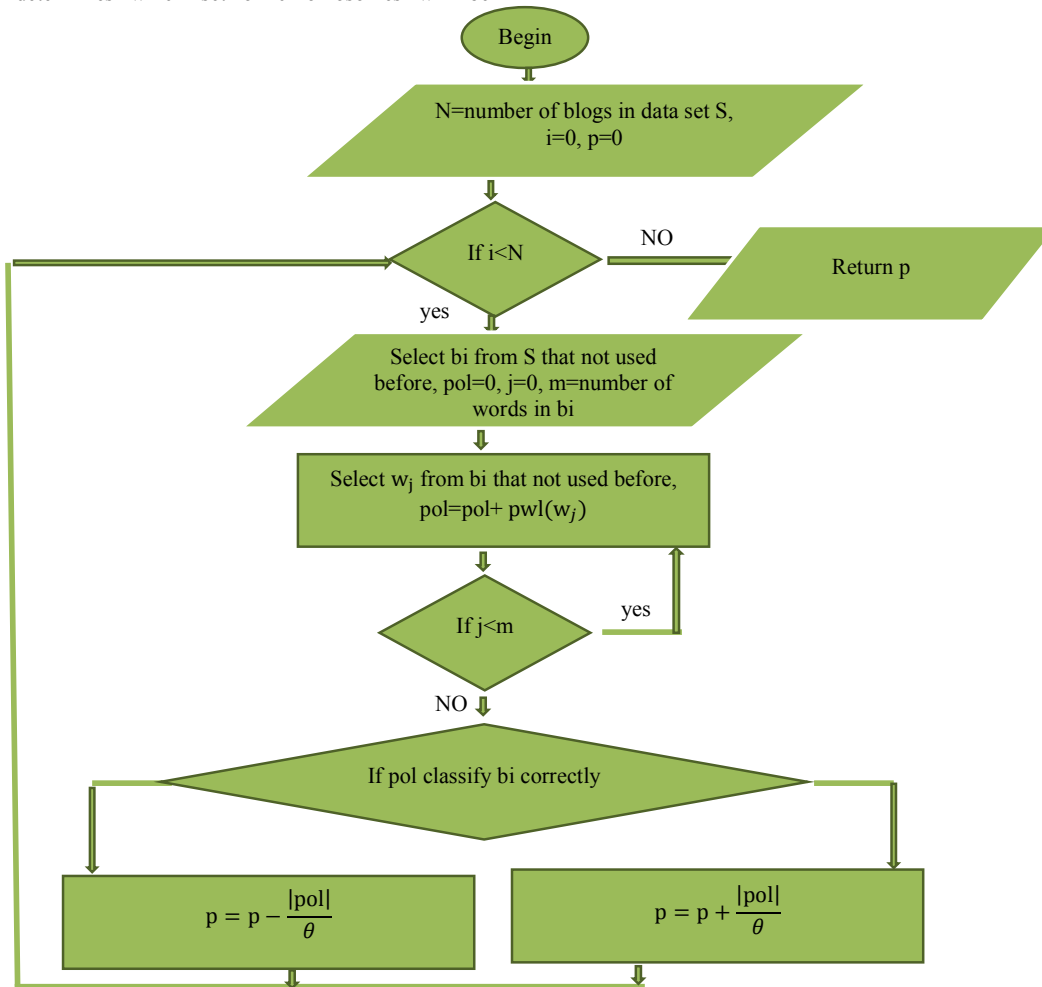


Fig 2 Fitness function calculation

3. POWER LINKS

Power links different from the ordinary frequency, instead of counting the frequency of individuals words, one counts the cooccurrence of two words. Power links is simpler the 2-gram counting, since the algorithm calculates the Power Link based on the appearance of the words in the same document, sentence or micro-blog. The relative distance between the words in the micro-blogs might affect the relation between the words, however, to reduce the calculation overload and since the length of the blogs is short, the Power Link is independent from the distance between words. The distance between words will not considered in the calculation.

In this paper, the Power Link is defined for two different effective words as the number of blogs that the two words appeared in divided by the total number of blogs. Effective words mean non-stop list words.

Let w_1 and w_2 are non-stop list words, where w_1 and w_2 appears simultaneously in $N_{w_1w_2}$ blogs and N is the total number of blogs in a given corpus, then the Power Link between w_1 and w_2 ,

$$PL(w_1, w_2) = \begin{cases} \frac{N_{w_1w_2}}{N} & \text{if } w_1, w_2 \text{ polarity of the same type} \\ 0 & \text{if } w_1, w_2 \text{ has different polarity} \end{cases}$$

Since, the lengths of the blogs are similar, no need to normalize the Power Link based on the lengths of the blogs. For each word w the algorithm calculate the Power Link between w and other words in the current corpus. This will form a matrix M of dimension $N \times N$, where N is the total number of words.

The lexicon that will be used will be produced in the first phase of the algorithm.

3.1 Classifying based on power links.

Classifying of blogs based on the Power Links is a clustering algorithm. The algorithm begins by selecting one of the blogs then calculates the distance between this blog and each blog in both classes $R1+$ and $R1-$. The blog will be added to the class with minimum average distance. This process will continue till all blogs are exhausted.

3.2 Distance between two blogs

For each blog b , define a vector blog vb , the length of vb is N . The element k in vb is the average of the corresponding row in M if $w_k \in b$, else the value is 0. The distance between two blogs b_1, b_2 is the Euclidian distance between the corresponding blog vectors vb_1, vb_2 . Fig4 illustrates the details of clustering algorithm based on Power Links.

4. EXPERIMENTS.

4.1 Data sets

In the experiments, a diversity data sets and multilingual is used. Five sets in English. These sets are Healthcare reform (HCR), SemEval dataset, Sanders –Twitter sentiment corpus, The Stanford Twitter dataset (STS dataset) and Obama-McCain debate (OMD) dataset [8] and another five data sets in Arabic. These sets are the crises of Turkish Lira (TLC), Muslims brotherhood (MBH), New High School regulations in Egypt (NHS), Egyptian elections (SIE) and American

elections (TRE) [7]. Table 1 shows the distribution of each set for a negative and positive blog.

4.2 Experimental Steps.

In our experiments, the value of minpol and maxpol are -10 and 10 respectively. The crossover and mutation are done based on 0.01 probability. The stop words list is calculated based on the frequency of each word. If a word appeared in all data sets and in each data set this word appeared in more than the half of the blogs then this word will be considered as a stop list word

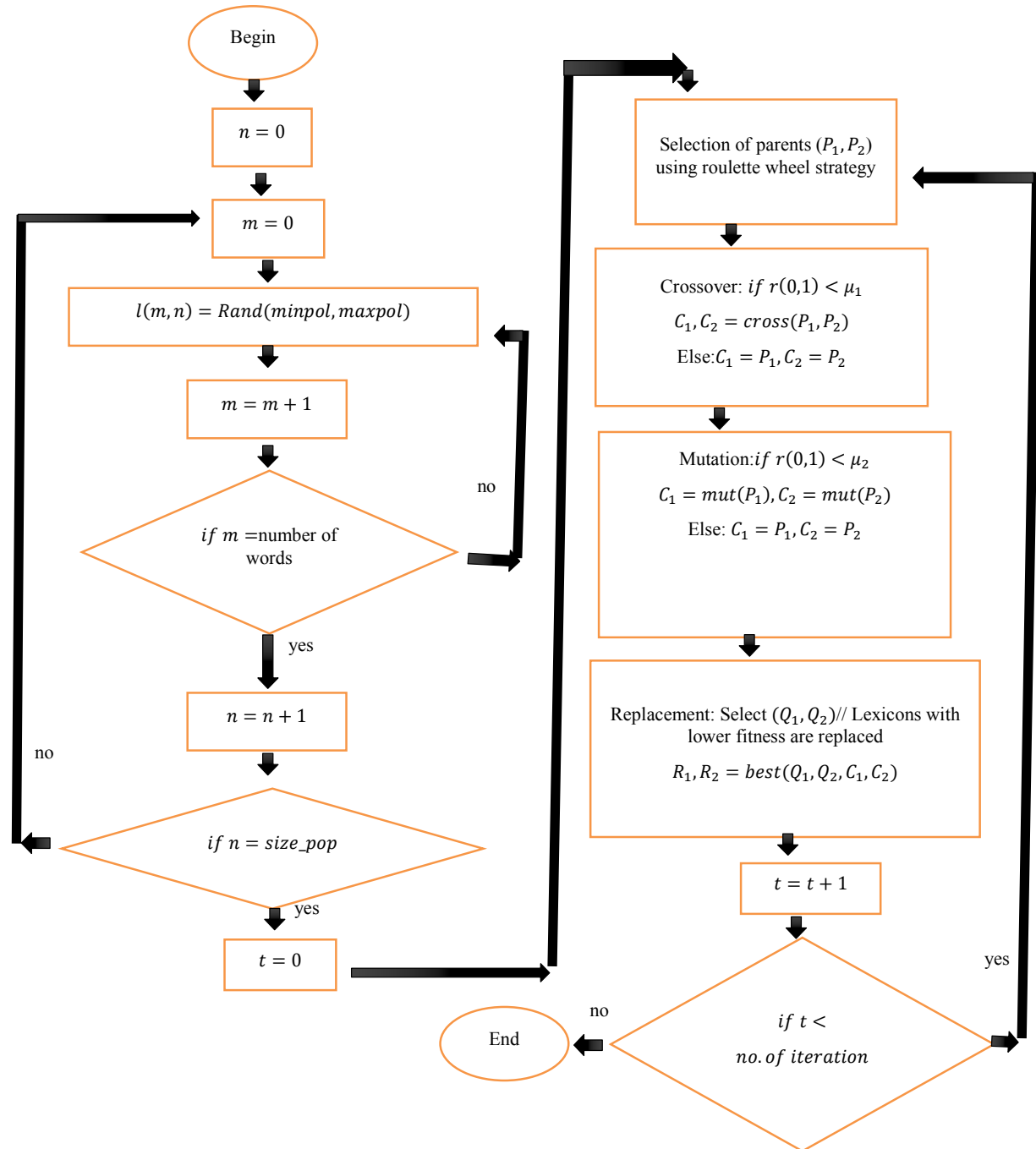


Fig.3 Genetic algorithm to optimize the accuracy of the lexicon

Table1. Negative, positive and total number of blogs in the data sets

	+ve	-ve	Total
HCR	917	369	1286
SemEval	3640	1458	5098
STS	182	177	359
Sanders	570	654	1224
OMD	710	1196	1906
TLC	431	271	702
MBH	631	442	1073
NHS	387	226	613
SIE	157	451	608
TRE	310	672	982

The experiments go according to the following steps:

1. One of the data sets is chosen randomly and the blogs in this data set are classified based on Bing-Lui lexicon.
2. The genetic algorithm is used to produce the next lexicon
3. The Power Link clustering algorithm is used to redistribute the blogs
4. The genetic algorithm is used again to produce a refined lexicon
5. The current lexicon is used to classify a new data set
6. The genetic algorithm will be used to produce a refined lexicon with one restriction, the algorithm will work to calculate the polarity of new words that never appear in the previous calculations
7. The steps 3, 4 and 6 are repeated to classify and refine the lexicon till all data sets are exhausted
8. The steps from 1 to 8 are iterated and in each iteration, begin with a different data set. In each iteration, the precision, recall and F-measure are calculated.
9. The iteration that gave the best precision, recall and F-measure is chosen to be the standard lexicon.

The results of these experiments will be tested based on repeating the experiments according to the following:

- A. Performing a classifying based on Bing-Lui lexicon only (BLC) [8]
- B. Performing a classifying based on random search method (RS) [8]
- C. Performing steps from 1 to 9 without using the Power Link clustering (GLO)
- D. Refining the results of 1 using the Power Link clustering (BLCP)
- E. Performing the experiments based on steps from 1 to 9 (GLOP)

Note that A and B are a supervised method. In these methods it is needed to assign a set for training and a set for testing. In these experiments, a 10-fold method is used. The whole data sets are divided randomly into 10 parts. In each fold, one of the folds is used as a test set and the others are used as a training set. The results below report for the average results for each data sets across the 10 testing iterations.

5. RESULTS AND DISCUSSION

In this section, results for the experiments sets A to E are reported.

Tables 2 to 11 reports the results of the proposed 5 methods for the 10 data sets. Also, the genetic algorithm was implemented using a different value of crossover and mutation probability, μ_1 and μ_2 respectively.

Table 2. Recall, Precision and F values of the data set HCR

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	45.5 7%	45.7 7%	45.6 7%	75.9 4%	77.4 0%	76.6 6%	61.1 7%
RS	50.8 8%	50.1 1%	50.4 9%	65.3 6%	64.9 6%	65.1 6%	57.8 3%
GLO	94.9 5%	93.6 7%	94.3 1%	85.6 4%	86.2 2%	85.9 3%	90.1 2%
BLC P	54.3 1%	51.7 4%	52.9 9%	81.8 0%	79.4 0%	80.5 8%	66.7 9%
GLO P	66.4 0%	67.1 2%	66.7 6%	75.0 3%	74.4 0%	74.7 1%	70.7 4%

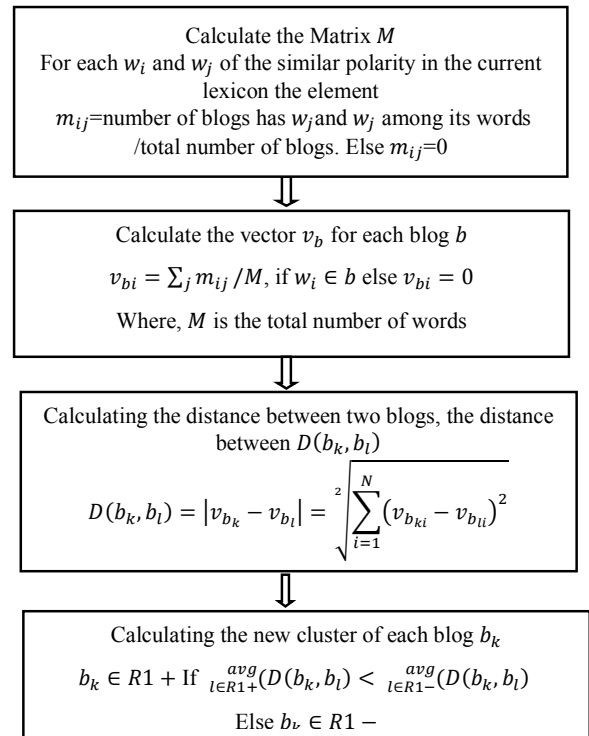


Fig4. Power Link clustering steps

The bold values show the values of the method that outperformed the other methods for the same data set based on F-measure. The results show that GLOP method outperformed all other methods in eight of the data sets. The values of μ_1 and μ_2 affects the convergence rate but these values should be less than 0.1 else the algorithm will never converge.

Table 3. Recall, Precision and F values of the data set SemEval

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	71.0 6%	71.9 1%	71.4 8%	40.2 6%	42.3 9%	41.3 0%	56.3 9%
RS	50.4 2%	48.9 7%	49.6 8%	38.7 1%	39.7 4%	39.2 2%	44.4 5%
GLO	72.4 2%	74.2 0%	73.3 0%	87.5 4%	85.2 0%	86.3 5%	79.8 3%
BLC P	69.6 3%	70.9 5%	70.2 8%	66.5 4%	68.7 0%	67.6 0%	68.9 4%
GLO P	83.4 6%	81.0 3%	82.2 3%	86.0 9%	85.6 4%	85.8 6%	84.0 5%

The number of iterations for convergence was between 150000 iterations and 250000 iterations.

Table 4. Recall, Precision and F values of the data set STS

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	84.8 3%	82.6 3%	83.7 2%	68.4 1%	66.2 9%	67.3 3%	75.5 2%
RS	66.6 1%	66.8 7%	66.7 4%	71.5 7%	71.3 7%	71.4 7%	69.1 0%
GLO	74.2 5%	74.0 7%	74.1 6%	55.0 3%	54.1 0%	54.5 6%	64.3 6%
BLC P	71.4 4%	69.9 1%	70.6 7%	87.6 9%	84.8 8%	86.2 6%	78.4 6%
GLO P	95.0 3%	93.4 6%	94.2 4%	87.6 6%	89.1 1%	88.3 8%	91.3 1%

It is noted that the number of iterations increase when the value of μ_1 and μ_2 decrease. Also, it is noted that merging the Power Link clustering method improved the results of the methods in average.

The main results that can be reported here is that a non-supervised method could present a performance that equal to or even overcome the performance of a supervised methods.

It is easily to note that the results of English data sets were better than the results of Arabic data sets. This is due to the quality of the initial used lexicons in each case and also due to the difficulty in classification of many Arabic blogs to a

correct +ve or -ve class manually. Many of the blogs gave no clear sentiment or gave +ve and -ve sentiment in the same time.

Table 5. Recall, Precision and F values of the data set Sanders

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	50.4 4%	48.3 8%	49.3 9%	35.0 0%	36.6 3%	35.8 0%	42.5 9%
RS	77.6 7%	75.4 4%	76.5 4%	57.8 5%	56.9 9%	57.4 2%	66.9 8%
GLO	76.3 8%	77.6 8%	77.0 2%	56.4 9%	57.2 6%	56.8 7%	66.9 5%
BLC P	58.1 1%	56.9 3%	57.5 1%	65.7 0%	65.1 6%	65.4 3%	61.4 7%
GLO P	88.1 6%	89.6 2%	88.8 8%	92.8 0%	95.5 5%	94.1 5%	91.5 1%

Table 6. Recall, Precision and F values of the data set OMD

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	57.1 2%	58.0 0%	57.5 6%	38.7 1%	87.0 8%	53.6 0%	55.5 8%
RS	54.8 7%	53.0 3%	53.9 3%	38.2 8%	60.6 0%	46.9 2%	50.4 3%
GLO	86.5 3%	87.4 9%	87.0 1%	61.7 7%	84.3 5%	71.3 2%	79.1 6%
BLC P	82.5 7%	85.1 0%	83.8 2%	46.8 0%	87.6 8%	61.0 3%	72.4 2%
GLO P	77.6 9%	78.7 6%	78.2 2%	55.4 8%	82.2 6%	66.2 7%	72.2 4%

Table 7 Recall, Precision and F values of the data set TLC

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	74.7 5%	77.6 4%	76.1 7%	39.2 1%	36.2 2%	37.6 6%	56.9 1%
RS	74.2 4%	72.9 3%	73.5 8%	45.0 2%	43.3 9%	44.1 9%	58.8 8%
GLO	83.0 5%	84.2 5%	83.6 5%	55.6 3%	57.3 8%	56.4 9%	70.7 1%
BLC P	69.1 1%	67.2 5%	68.1 7%	72.4 1%	72.0 9%	72.2 5%	70.2 1%
GLO P	85.2 0%	84.6 5%	84.9 2%	79.1 7%	77.0 7%	78.1 1%	80.8 8%

Table 8 Recall, Precision and F values of the data set MBH

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	75.2 6%	75.3 9%	75.3 2%	43.4 0%	41.1 6%	42.2 5%	58.7 9%
RS	72.8 5%	75.6 5%	74.2 2%	42.6 0%	45.2 8%	43.9 0%	59.0 6%
GLO	79.2 2%	80.3 7%	79.7 9%	53.9 7%	53.8 9%	53.9 3%	66.8 6%
BLC P	50.8 2%	52.9 1%	51.8 4%	58.7 2%	59.7 8%	59.2 5%	55.5 4%
GLO P	81.8 0%	80.2 6%	81.0 2%	75.3 4%	76.5 5%	75.9 4%	78.4 8%

Table 9 Recall, Precision and F values of the data set NHS

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	57.5 5%	58.8 2%	58.1 8%	42.6 0%	43.3 2%	42.9 6%	50.5 7%
RS	73.3 0%	70.8 2%	72.0 4%	59.2 2%	57.1 3%	58.1 6%	65.1 0%
GLO	53.4 6%	56.2 7%	54.8 3%	61.3 4%	62.1 5%	61.7 4%	58.2 9%
BLC P	79.3 0%	77.2 4%	78.2 6%	77.0 3%	77.0 6%	77.0 4%	77.6 5%
GLO P	79.6 2%	79.5 2%	79.5 7%	71.0 2%	70.3 1%	70.6 6%	75.1 2%

Table 10 Recall, Precision and F values of the data set SIE

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	50.0 2%	51.6 3%	50.8 1%	62.6 8%	63.4 5%	63.0 6%	56.9 4%
RS	41.4 3%	41.4 6%	41.4 4%	67.4 8%	68.3 5%	67.9 1%	54.6 8%
GLO	77.0 8%	79.7 6%	78.4 0%	59.8 9%	58.5 5%	59.2 1%	68.8 0%
BLC P	69.7 2%	68.4 2%	69.0 6%	84.6 2%	84.1 2%	84.3 7%	76.7 2%
GLO P	82.3 9%	82.8 4%	82.6 1%	83.0 5%	84.0 9%	83.5 7%	83.0 9%

Table 11 Recall, Precision and F values of the data set TRE

Met hod	Positive Class			negative Class			AV G F
	R	P	F	R	P	F	
BLC	62.4 7%	60.1 9%	61.3 1%	58.6 9%	56.1 6%	57.4 0%	59.3 5%
RS	70.6 0%	72.5 5%	71.5 6%	36.8 3%	38.9 7%	37.8 7%	54.7 2%
GLO	67.3 1%	65.7 5%	66.5 2%	59.6 7%	61.7 3%	60.6 8%	63.6 0%
BLC P	54.6 0%	52.8 0%	53.6 8%	50.7 8%	51.0 9%	50.9 3%	52.3 1%
GLO P	89.3 9%	89.6 6%	89.5 2%	91.6 7%	93.8 6%	92.7 5%	91.1 3%

The quality of the classification of the proposed method is affected by the initial lexicon that should be prepared in advance, however regardless the quality of this lexicon the algorithm was able to present an accepted results even with using a random lexicon at the beginning.

6. CONCLUSION

In this work unsupervised method to determine the polarity of blogs were presented. The method overcome the performance of other methods either supervised or unsupervised methods in most of the data sets. The performance of the algorithm is affected by the initial lexicon and the language of the blogs. The future work will concentrate on developing the algorithm to work without the need of initial lexicon. Also the algorithm will be expanded to increase the number of polarity classes.

7. REFERENCES

- [1] Sarangi, S. K., Jaglan, V. and Dash, Y., 2013. A Review of Clustering and Classification Techniques in Data Mining, International Journal of Engineering, Business and Enterprise Applications (IJEBA), 4(2), March-May, , pp. 140-145
- [2] Aucouturier, J.-J. and Pachet, F., 2004, "Improving Timbre Similarity: How high's the sky?," in Journal of Negative Research Results in Speech and Audio Sciences.
- [3] Kumar, V., and Rathee, N., 2011, ITM University, "Knowledge discovery from database Using an integration of clustering and classification", International Journal of Advanced Computer Science and Applications, (March 2011) Vol. 2, No.3.
- [4] Sharma, N., Bajpai, A. and Litoriya R., 2012, "Comparison the various clustering algorithms of weka tools", International Journal of Emerging Technology and Advanced Engineering, ISSN 2250-2459, (May 2012), Volume 2, Issue 5.
- [5] Devijver, P.A. and Kittler, J., 1982, Pattern Recognition: A Statistical Approach. PrenticeHall.
- [6] H. Liu and H. Motoda, 1998, Feature Extraction, Construction and Selection: A Data Mining Perspective. Kluwer Academic.

- [7] Rokaya, M., Ghiduk A. S., 2019, Arabic Lexicon Learning to Analyze Sentiment in Microblogs, (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 10, No. 8, 592-599
- [8] Keshavarz, H. and Abadeh, M. S., 2017, ALGA: Adaptive lexicon learning using genetic algorithm for sentiment analysis of microblogs, Knowledge-Based Systems 122 , 1–16
- [9] Haj mohammadi, M.S., Ibrahim and Selamat R. A., 2014, Cross-lingual sentiment classification using multiple source languages in multi-view semi-supervised learning, Eng. Appl. Artif. Intell. 36, 195–203 .
- [10] Cui, Z., Shi, X. and Chen, Y., 2015, Sentiment analysis via integrating distributed representations of variable-length word sequence, Neurocomputing, 126–132 .
- [11] Weissbock, J. and Inkpen, D., 2014, in: Combining Textual Pre-Game Reports and Statistical Data for Predicting Success in the National Hockey League, Advances in Artificial Intelligence, Springer International Publishing, pp. 251–262 .
- [12] Javed, B.S., 2018, Hybrid semantic clustering of hashtags, Online Social Networks and Media 5 (2018) 23–36
- [13] Vicient , A. M., 2014, Unsupervised semantic clustering of Twitter hashtags, Proceedings of the 21st European Conference on Artificial Intelligence, pp. 1119–1120 .
- [14] Javed , B.S., 2016, Sense-level semantic clustering of hashtags in social media, in: Proceedings of the 3rd Annual International Symposium on Information Management and Big Data, pp. 140–149 .
- [15] Muntean, C.I., Morar , G.A., Moldovan, D., 2012, Exploring the meaning behind Twitter hashtags through clustering, Lect. Notes Bus. Inf. Process. 127, 231–242
- [16] Bhulai, S., Kampstra, P., Kooiman, L., Koole, Deurloo, G., M. and CCing, B.K., 2012, Trend visualization on Twitter: what's hot and what's not?, Proceedings of the 1st International Conference on Data Analytics, Springer-Verlag, pp. 43–48
- [17] Goldberg, D. E., 1989, Genetic Algorithms in Search, Optimization, and Machine Learning. Reading: Addison-Wesley
- [18] Pesina, S. and Yusupova, L. G., 2014, Words Functioning in Lexicon, 2nd Global Conference On Linguistics And Foreign Language Teaching, Linelt-, Dubai – United Arab Emirates, December 11 – 13, 2014
- [19] Rokaya, M and Atlam, E. S., 2010, Building of field association terms based on links', Int. J. Computer Applications in Technology, Vol. 38, No. 4, pp.298–305.
- [20] Rokaya, M, 2013, Automatic Text Extraction Based on Field Association Terms and Power Links, International Journal of Computer and Information Technology (IJCIT), Volume 02– Issue 06 (November 2013), pp 1049- 1053
- [21] Rokaya, M, 2013, Automatic Summarization based on Field Coherent Passages. International Journal of Computer Applications 79(9) (October 2013),38-44., Published by Foundation of Computer Science, New York, USA
- [22] Rokaya, M and Aljahdali, S., 2013, Building a Real Word Spell Checker Based on Power Links, International Journal of Computer Applications, Vol. 65 No 7,(March 2013), PP 14-19.
- [23] Rokaya, M, Nahla, A. and Aljahdali, S., 2012, Context-Sensitive Spell Checking Based on Field Association Terms. IJCSNS International Journal Of Computer Science And Network Security. Vol. 12 No. 3 pp. 64-68.
- [24] Rokaya, M, 2014, Improving Ranking of Search Engines Results Based on Power Links, IPASJ International Journal of Information Technology (IJIT),Volume 2, Issue 9, September.