# Hybrid Proposed Method Using Statistical Linguistic Features (SLF) And Neural Network In Arabic Texts Summarization

*Iman Qays Abduljalee*      *Amal Hameed Khaleel*      *suhad muhajer kareem*

*Basra university \ Science collage \Computer science dept.*

*emankais@yahoo.com*

## Abstract

With the counting growth of electronic information on World Wide Web, it has become necessary important to provide mechanisms to find and present a shorter version would suffice, so that automatic text summarization technique which plays an important role to help users to determine whether it has to do with information they need or not.

In this paper we present a short historical overview and advancement of automatic text summarization and the most relevant approaches currently used in this area. We proposed a new technique in automatic summarization area for Arabic news articles using a neural network by selecting important sentences from the original text and put it in the summary. A Multi-layer Perceptron neural network (MLP) is trained to learn the relevant characteristics of sentences that should be included in the summary of the article. The neural network is then modified to generalize and combine the relevant characteristics apparent in summary sentences . Finally, system is evaluated by comparing the final summary of the system with the summary produced by expert in Arabic language, we measure the performance of the system by computing the precision, recall and F-Measure and we obtain good result that we display in the conclusion.

**Keywords:** Text summarization, Multi-layer Perceptron Neural Network (MLPNN), F-Measure.

## I.    Introduction

Text summarization has become an important and timely tool for assisting and interpreting text information in today's fast-growing information age. It is very difficult for human beings to manually summarize large docu-ments of text. There is an abundance of text material available on the internet(V. Gupta *et al*.2010). Text summarization is the process to produce a condensed representation of the content of its input for human consumption (S.D. Afantenos *et al*.2005).

   The input to a summarization system can be one or more text documents, When only one document is the input, it is called single document text summarize-ation and when the input is a cluster of related text documents, it is called multidocument summary-zation. We can also categorize the text summarization based on the type of users the summary is

intended for query focused summaries are tailored to the requirements of a particular user or group users and generic summaries are aimed at a broad readership community (K. Sarkar 2009). Text Summarization methods can be classified into extractive and abstractive summarization. An extractive summarization method consists of selecting important sentences, paragraphs etc. from the original document and concate-nating them into shorter form. The importance of sentences is decided based on statistical and linguistic features of sentences. An Abstractive summarization attem-pts to develop an understanding of the main concepts in a document and then express those concepts in clear natural language. It uses linguistic methods to examine and interpret the text and then to find the new concepts and expressions to best describe it by generating a

new shorter text that conveys the most important information from the original text document. This paper focuses on extractive text summarization methods. Extra-ctive summaries are formulated by extracting key text segments (sentences or passages) from the text, based on statistical analysis of individual or mixed surface level features such as word/phrase frequency, location or cue words to locate the sentences to be extracted. The "most important" content is treated as the "most frequent" or the "most favorably positioned" content. Such an approach thus avoids any efforts on deep text understanding. They are concept-ually simple, easy to implement (V. Gupta *et al*.2010, J. Lin 2009, K. Jezek and J. Stteinberger 2007).

## II.     Related Works

Most early work on single summarization system focused on frequency word proposed by luhn in 1958(H. P. Luhn 1958), his research based on high frequency to find important parts from original document to generate summary. In 1969, where edmunson used cue words and the similarity to the title as two features to extract salient parts from text. In 2001 lin & gon used latent semantic analysis to extract the meaning from the words and the sentences by using model of singular vector decomposition (SVD) (J. Steinberger 2001). A new method used in text summarization using graph-based method proposed by Radav & Mihalca in 2004 to present the text in graph architecture where nodes are the sentence and the edges  are the similarity between sentences (R. Mihalcea 2004).
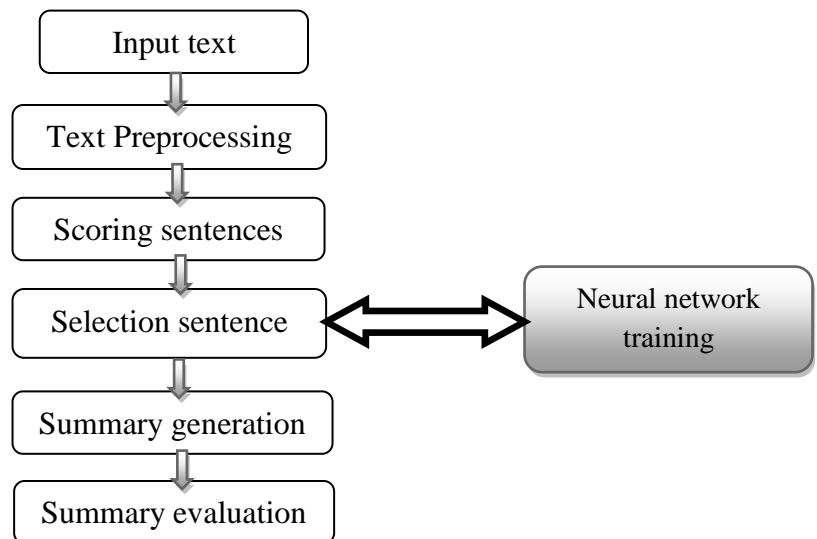
After that a numbers of learning techniques are proposed

and used in text summarization such as: kupiek in 1995 used Bayesian classifier to compute the probability of each sentence in the text to appeared in the summary(J. Kupiek *et al*. 1995), Khosrow kaikhah in 2004 used neural networks to train and learn the relevant features of sentences that should be included in the summary of text documents (Khosrow 2004),in 2011 Ben King *et al.*  proposed extractive text summarization using neural network(Ben King *et al.* 2011),

and numbers of researches are still proposed in text summarization.

## III.   The   Proposed   Method Architecture

In this papers, we proposed a new method to generate summary for article Arabic texts by using neural network where the method consist of numbers of stages to generate summary and we shown in the following diagram figure (1):-

```
          ┌─────────────┐
          │ Input text  │
          └─────────────┘
                 ↓
        ┌──────────────────┐
        │ Text Preprocessing│
        └──────────────────┘
                 ↓
        ┌──────────────────┐
        │ Scoring sentences │
        └──────────────────┘
                 ↓
        ┌──────────────────┐        ┌──────────────────┐
        │ Selection sentence│ ⟺      │ Neural network   │
        └──────────────────┘        │ training         │
                 ↓                   └──────────────────┘
        ┌──────────────────┐
        │ Summary generation│
        └──────────────────┘
                 ↓
        ┌──────────────────┐
        │ Summary evaluation│
        └──────────────────┘
```

Figure(1)  the proposed method architecture .

The model consist of the following stages:-

**1-** Text Preprocessing stage:- this stage consist of :

**1-1** tokenization:- process to determines the boundaries of each word and sentence in the original text.

**1-2** stemming:-

1-3 we obtain the stem of each word by removing suffixes from this word such as (درس ، يدرس ←درس ، ذهب ←يذهبون ….est. )

1-4 part of speech (word tagging):- is the process of assigning POS to each word in the sentence to give the class of this word such as (noun, verb, adjective,…ect). This process is depend on the lexical stored in the system.

1-5 stop words removal:- is the process to filter the text from the most frequent important words by using list of stop word stored in the system such as ( التي ، ......... ، الذي، ماذا، الجدير بالذكرect). The advantage of removing this word to

increasing the performance of computing scoring of the sentence.

**2-** Sentences Scoring stage:-is the process to compute the weight of each sentence based on number of features we are called as SLF(Statistical Linguistic Features) features of texts, each sentence is presented as vector $[F_1, F_2, F_3, F_4, F_5, F_6, F_7]$, where:-

$F_1$ *(frequency word feature) :- compute number of the occurrences of words in full original text.*

$$F1 = \sum_{i=1}^{n} frequency(w)$$

$F_2$ *(length sentence feature):- compute number of words within its sentence without stop of word.*

$$F2 = \sum_{i=1}^{n} no.\,of\,words(s)$$

$F_3$ *(position sentence feature):- the location sentence in the text.*

$$F3 = \frac{1}{position(s)}$$

$F_4$ *(no. Verbs feature):- compute the number of verbs in the sentence of input text based on*

*POS.*

$$F2 = \sum_{i=1}^{n} no. \, of \, verbs(s)$$

**$F_5$ (no. proper names feature):-** *compute the number of proper names in the sentence of input text based on POS.*

$$F2 = \sum_{i=1}^{n} no. \, of \, proper \, names(s)$$

**$F_6$ (no. digits feature):-** *compute the number of digits in the sentence of input text.*

$$F2 = \sum_{i=1}^{n} no. \, of \, digits(s)$$

**$F_7$ (similarity to title feature):-** *compute the number of common words between the word of sentence and the word of title of input text.* *F7==Σ SIM(Si , Sj) i>=j*

**3-** Sentences Selection stage:- we used neural network method that  is performed in two different phases namely, training phase and testing phase.  The first phase of the process involves training the neural networks to learn the types of sentences that should be included in the summary. This is accomplished by training the network with sentences where each sentence is identified as to whether it should be included in the summary or not. This is done by a human reader. We use a three-layered    feedforward MLP neural   network.   Our   neural network consists of seven input-layer neurons, five hidden-layer neurons,   and   one   output-layer neuron ( as shown in figure 2). We uses   the   sigmoid   activation function for both the weights of the hidden layer and the output layer.

*Algorithm*

i.    Create an architecture consists of three layers: input, hidden, and output layers, for input layer, assign as net input to each unit ($x_i$, i=1,…,7) its corresponding features in the input vector. The output for each unit is its net input.

ii.   Initialize the weights and bias to random values .

iii.   Initialize the network parameters. ( Learning rate =0.3, Momentum rate =0.7)

iv.   Train the network with initialized parameters, and with sigmoid activation function.

v.   Calculate the error using MSE method

vi.  Repeat the process until the maximum epochs are reached or the desired output is identified or the minimum gradient is reached.  (Maximum epochs = 10000, Total error = $10^{-6}$)



Figure(2) :-Structure of Multi-layer Perceptron Neural Network

**4-**    Summary generation stage:- In this step, after the scores of sentences were selected by neural network to appear in the summary, we use the following steps to generate the summary:-

*Summary length =  Summarization Rate*
*\*(Document Length) /100  ….(1)*

- Sort this scores of selected sentences by neural network in decreasing order  and select n top of this score, the value of n depend on summary length which computed from the following formula:-

  In this paper, we used summarization rate is 45%.

- Sentences may not be ordered in the correct order so we should ranked it in order as the order appeared in the original texts and we put title of the original text as the title of the summary. The advantage of this

step to make summary more readable and coherent.

**5-**    Summary evaluation stage:- summaries can be evaluated using intrinsic or extrinsic measures, while intrinsic methods attempt to measure summary quality using human evaluation therefore, extrinsic methods measure the same through a task-based measure such the information retrieval-oriented task(D. Mallett, J. Elding and Mario 2004).

In this work we used ROUGE (Recall-Oriented Understudy Gisting Evaluation ) measure is one of most intrinsic methods  to evaluate summary which written by human(expert human) with automatic summary to measure performance of the summary. This measure produce numerical value by use precision(P),recall(R) and F-Measure(F)( R. M. ALGuliev and R. M. ALGullyev 2007):-

$$P = \frac{Sman \cap Sauto}{Sauto} \qquad \ldots(2)$$

$$R = \frac{Sman \cap Sauto}{Sman} \qquad \ldots(3)$$

$$F\text{-}Measure = \frac{2 \times P \times R}{P + R} \qquad \ldots(4)$$

Where  :- **Sman** is the manual summary and **Sauto** is the automatically-generated summary.

### Examples :-the input text to summarize it.



Figure (3) : the input Arabic text

## IV.    Experiment Result

Where each document is converted into a list of sentences. Each sentence is represented as a vector $[f_1, f_2, ..., f_7]$, composed of 7 features.

The selection of features plays an important role in determining the type of sentences that will be selected as part of the summary and, therefore, would influence the performance of the neural network. The proposed system is implemented in Delphi 7 . In this section we presented analysis and evaluation result.

## a.      Analysis results

We used 150 news articles from the Internet (i.e. Wikipedia) articles with various topics such as technology, sports, and world news to train the network. Each article consists of  20 to 100 sentences. Every sentence is labeled as either a summary sentence or an unimportant sentence by a human reader. We then used the same 50 news articles as a test set for the modified network. The accuracy of the modified network ranged from 91% to 100% with an average accuracy of 94.3% when compared to the summaries of the human reader, the modified network included a sentence that was not selected by the human reader. That is, the network was able to select all sentences that were labeled as summary sentence in most of the articles.

Figure 4 show the features scores for the input text  which it interred to system from (figure3) while figure 5 show feature score that choose from features scores for the text document using combines of statistical and linguistic features to compute final scores of each sentence in input text, this final scores presented to neural network to learn it and select sentences which select to generate final summary which show in figure 6.

Figure (4) : Features scores for the sentences of input Arabic text using combines of statistical and linguistic (SAL) features.

Figure  (5) :Feature score that choose from features scores using MLP neural network  for the input Arabic text
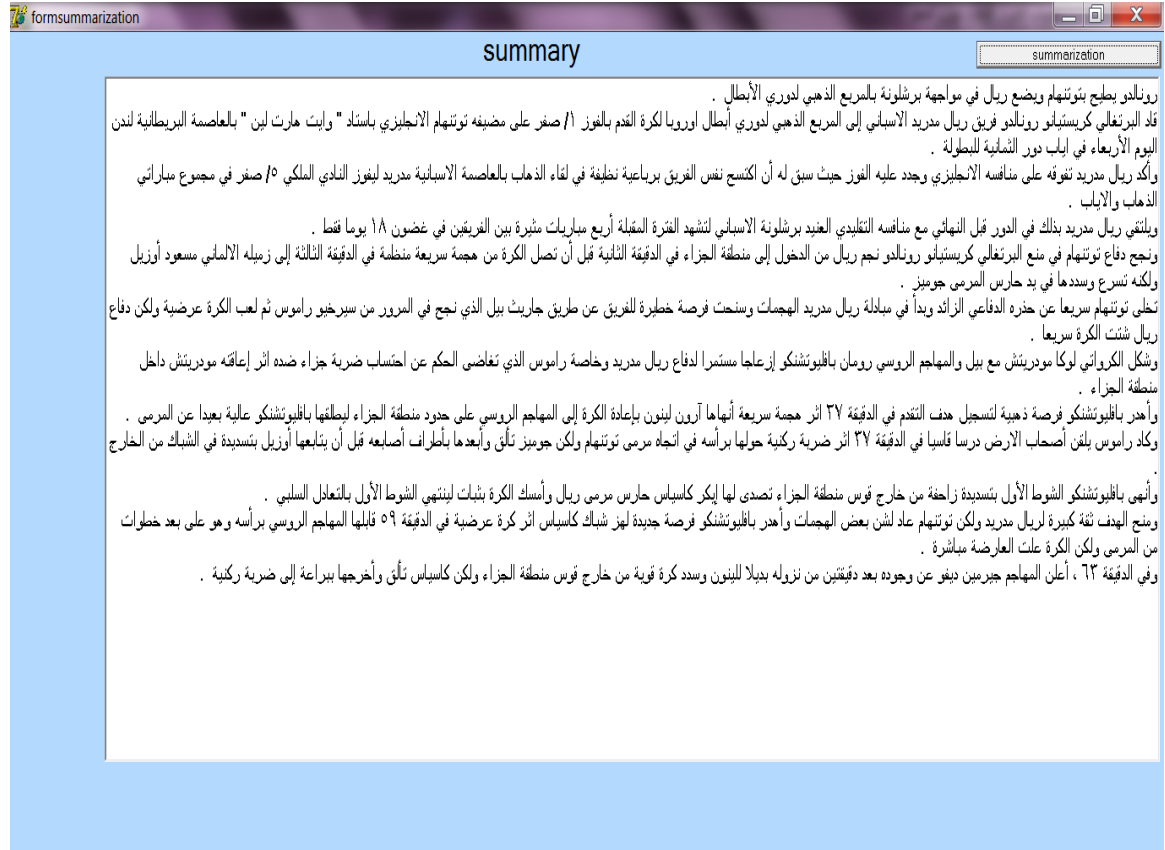


Figure (6):- summary of input Arabic text

### b.  Evaluation results

The performance of the proposed approach is evaluated using precision, recall and F-measure from formulas (2,3,4). Precision evaluates the proportion of correct of the sentences in the summary whereas recall is used to evaluate the proportion of relevant sentences included in summary. For precision, the higher the values, the better the system is in
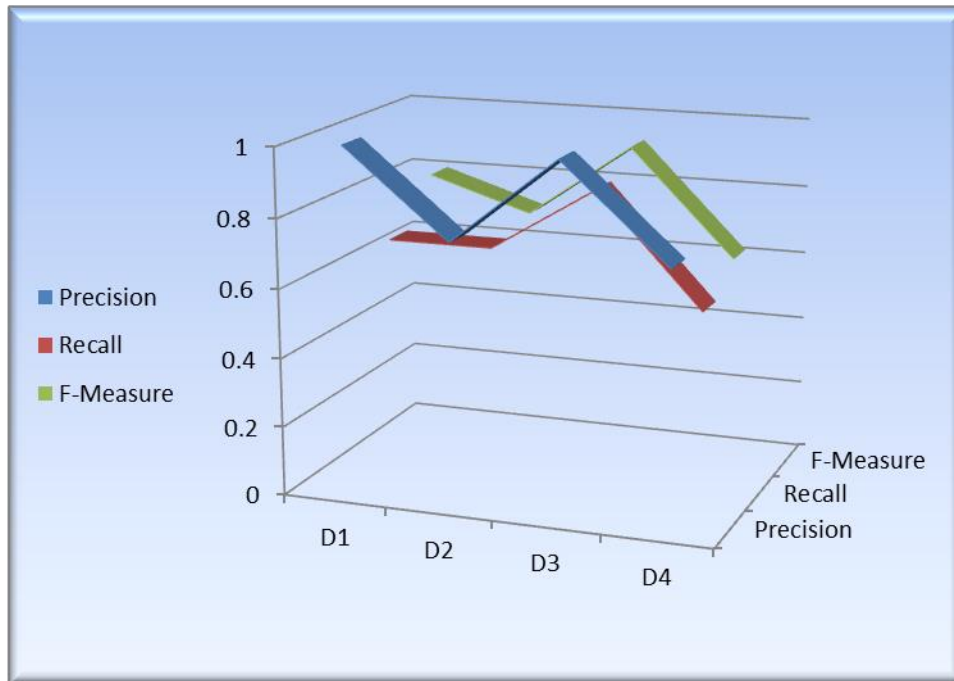
omitting irrelevant sentences and always precision values is higher than recall values. After computes precision and recall, we compute F-measure.

The table 3 show samples of the values for the evaluation measures(precision, recall and F-Measure) with summarization rate=45. In figure 7 we shown the diagram of precision, recall and F-measure values. In the table4 ,we display two results of other researches

Table 3: the values for the evaluation measures for four document

| Documents | Precision | Recall | F-Measure |
|-----------|-----------|--------|-----------|
| $D_1$ | 1 | 0.666 | 0.8 |
| $D_2$ | 0.75 | 0.666 | 0.705 |
| $D_3$ | 1 | 0.857 | 0.923 |
| $D_4$ | 0.727 | 0.533 | 0.615 |

Figure(7) :-diagram the values of precision, recall and F-measure

Table 4: the results of numbers of  evaluation measures for other research

| Research no. | authors | year | Performance of method |
|---|---|---|---|
| 1 | Khosrow Kaikhah | 2004 | 0.91 |
| 2 | Ben King   et.al | 2011 | 0.459 |

## V.    Conclusions

Our text summarization method performed well on the test article, with accuracy of 94%. The selection of features as well as the selection of summary sentences by the human reader from the training paragraphs plays an important role in the performance of the network. The network is trained according to the style of the human reader

and to which sentences the human reader deems to be important in a article. Individual readers can train the neural network according to their own style. In addition, the selected features can be modified to reflect the reader's needs and requirements. each property an important role in the selection of sentences where there is a relation between the properties of text, for example, when combining two or more properties with each other, the results were different appear. Where we found when combining property length with property similarity to the title the result was not good, as to when

combining the Length property and position, they were give better results than its the previous, so we tried hard to find all the properties that, when combined we get excellent results for the bottom and we found that these features when combined get a summary of a strong and coherent as it appears in the results of the evaluation. In the future we will try to develop the integrated system by combined the proposed neural network with genetic algorithms because the role of GAs in the best choice for the sentences that will appear in the final summary.

## References

**Ben King et al.**( 2011) "Experiments in Automatic Text Summarization Using Deep Neural Networks", In Proceeding of the Machine Learning, pp.1-12.

**D. Mallett, J. Elding and Mario,** (2004) "Information -Content Based Sentence Extraction for Text Summarization", In Proceeding of the IEEE, Canada, Vol.2(2), pp.1-5.

**H. P. Luhn**,(1958) "The Automatic Creation of Literature Abstracts", In Proceeding of IBM Journal of Research Development, Vol.2(2), pp.159-165.

**J. Kupiek, J. Pedersen and F.Chen**, .(1995) "Trainable Document Summa-rizer", In Proceeding of Research and Development in Information Retrieval, pp.68-73.

**J. Lin**,(2009) "Summarization", In Proceeding of Springer-Verlag, pp.1-7.

**J. Steinberger,** (2001)"Using Latent Semantic Analysis in Text Summarization and Summary Evaluation", In proceeding of the seventh International Conference ISIM, pp.1-8.

**K. Jezek and J. Stteinberger**, (2008)"Automatic Text Summar-ization (the State of Art 2007 and New Challenges)", In Proceeding

of Znalosti , pp.1-12.

**Khosrow Kaikhah** (2004), "Automatic Text Summarization with Neural Networks", In Proceeding on IEEE International Conference on Intelligent Systems, pp.40-44.

**K. Sarkar,** (2009)"Using Domain Knowledge for Text Summariza-tion In Medical Domain", In Proceeding of ACEEE, Interna-tional Journal of Recent Trends in Engineering, Vol.1, no.1, pp.200-205.

**R. Mihalcea**,(2004) "Language Independent Extractive Summa-rization ", In Proceeding of the ACL Interactive Poster and Development Sessions, pp.49-52.

**R. M. ALGuliev, R. M. ALGullyev** (2007)"Experimental Investigating the F-Measure as Similarity Measure for Automatic Text Summarization", In

Proceeding of the Computational Mathematical, Vol.6(2) , pp.278-287, 2007.

**S.D. Afantenos, V. Karkaletsis and P.Stamatopoulos**,(2005) "Summarization from Medical Documents: A Survey ", In Proceeding of *Artificial Intelligence in Medicine*, vol. 33, pp. 157-177.

**VI. Gupta and G. Lehal** ,(2010) "A Survey of Text Summarization Extractive Techniques ", In Proceeding of journal of emerging technologies in web intelligence, vol. 2, no. 3, pp. 258-268.

# خوارزمية هجينة مقترحة باستخدام الصفات الاحصائية اللغوية SLF والشبكات العصبية لتلخيص النصوص العربية

إيمان قيس عبد الجليل        أمل حميد خليل        سهاد مهجر كريم

جامعة البصرة \ كلية العلوم \ قسم علوم الحاسبات

emankais@yahoo.com

المستخلص

مع الزيادة المستمرة لنمو المعلومات الالكترونية المتوفرة على الشبكة العنكبوتية، أصبح من الضروري والمهم للحصول على آلية لتمثيل وتقديم نص بأسلوب مختصر ومفيد، لذلك ظهرت تقنية التلخيص الآلي للنصوص والتي تلعب دورا مهما في مساعدة المستخدم لتحديد ما يحتاجه من تلك النصوص.

في هذا البحث سنقدم نبذة تاريخية مختصرة عن تلخيص النصوص واغلب الطرق ذات الصلة المستخدمة حاليا في هذا المجال. حيث تم اقتراح تقنية جديدة في مجال التلخيص الآلي للمقالات الإخبارية العربية باستخدام مفهوم الشبكات العصبية. حيث تم اختيار الجمل المهمة من النص الأصلي ووضعها في الخلاصة ، ومن خلال تدريب شبكة Multi-layer Perceptron تم استخراج الخصائص ذات العلاقة من الجمل التي ينبغي أن تدرج في الخلاصة، وبعد ذلك تم تعديل الشبكة العصبية (MLP) بالتعميم والجمع بين الخصائص ذات الصلة لتظهر في جمل الخلاصة. وأخيرا، النظام تم تقييمه بمقارنة الخلاصة النهائية للنظام مع الخلاصة التي كتبها الخبير اللغوي في مجال اللغة العربية، وقمنا بقياس انجازيه النظام وذلك بحساب مقاييس الدقة و الاسترجاع ومقياس F-Measure. أثناء التقييم حصلنا على نتائج جيدة موضحة في الاستنتاج.