

من FCIH Wiki

- ١ Project data
- ٢ Team members
- ٣ Abstract
- ٤ Data set
- ٥ Methodology
 - ٥.١ Document representation Techniques Used:
 - ٥.٢ Classification Algorithms used :
 - ٥.٣ Experiment #1(CNN):
 - ٥.٤ Experiment #2(CNN):
 - ٥.٥ Experiment #3(CNN):
 - ٥.٦ Experiment #4 (CNN):
 - ٥.٧ Experiment #5(KNN):
 - ٥.٨ Experiment #6(KNN):
 - ٥.٩ Experiment #7(KNN):
 - ٥.١٠ Experiment #8(KNN):
 - ٥.١١ Experiment #9:
 - ٥.١٢ Experiment #10:
 - ٥.١٣ Experiment #11:
 - ٥.١٤ Experiment #12:
 - ٥.١٥ Experiment #13:
 - ٥.١٦ Experiment #14:
 - ٥.١٧ Experiment #15:
 - ٥.١٨ Experiment #16:
 - ٥.١٩ Experiment #17:
 - ٥.٢٠ Experiment #18:
 - ٥.٢١ Experiment #19:
 - ٥.٢٢ Experiment #20:
 - ٥.٢٣ Experiment #21:
 - ٥.٢٤ Experiment #22:
 - ٥.٢٥ Experiment #23:
 - ٥.٢٦ Experiment #24:
 - ٥.٢٧ Experiment #25:
 - ٥.٢٨ Experiment #26:
 - ٥.٢٩ Experiment #27:
 - ٥.٣٠ Experiment #28:
 - ٥.٣١ Experiment #29:
 - ٥.٣٢ Experiment #30:
 - ٥.٣٣ Experiment #31:
- ٦ Results
 - ٦.١ Experiment#1:
 - ٦.٢ Experiment#2:
 - ٦.٣ Experiment#3:
 - ٦.٤ Experiment#4:
 - ٦.٥ Experiment# 5 :
 - ٦.٦ Experiment# 6 :

- 7.7 Experiment# 7 :
- 7.8 Experiment# 8 :
- 7.9 Experiment #9:
- 7.10 Experiment #10:
- 7.11 Experiment #11:
- 7.12 Experiment #12:
- 7.13 Experiment #13:
- 7.14 Experiment #14:
- 7.15 Experiment #15:
- 7.16 Experiment #16:
- 7.17 Experiment #17:
- 7.18 Experiment #18:
- 7.19 Experiment #19:
- 7.20 Experiment #20:
- 7.21 Experiment #21:
- 7.22 Experiment #22:
- 7.23 Experiment #23:
- 7.24 Experiment #24:
- 7.25 Experiment #25:
- 7.26 Experiment #26:
- 7.27 Experiment #27:
- 7.28 Experiment #28:
- 7.29 Experiment #29:
- 7.30 Experiment #30:
- 7.31 Experiment #31:

http://fci.helwan.edu.eg/wiki/%D9%85%D8%B9%D8%A7%D9%84%D8%AC%D8%A9_%D8%A7%D9%84%D9%84%D8%BA%D8%A7%D8%AA_%D8%A7%D9

Our project based on text classification, Text classification is a core problem to many applications, like spam detection, sentiment analysis or smart replies. The goal of text classification is to assign documents (such as emails, posts, text messages, product reviews, etc...) to one or multiple categories. Such categories can be review scores, spam v.s. non-spam, or the language in which the document was typed. In our project we classify articles as we classify any input come to one category of (Business, science, health, entertainment) . our classification based on headline of article as we tokenize it from article and take the most effective term which can differentiate between article and another.

News Aggregator dataset contains headlines, URLs, and categories for 422,937 news stories collected by a web aggregator between March 10th, 2014 and August 10th, 2014.

News categories included in this dataset include business; science and technology; entertainment; and health. Different news articles that refer to the same news item (e.g., several articles about recently released employment statistics) are also categorized together.

Content The columns included in this dataset are:

ID : the numeric ID of the article TITLE : the headline of the article, URL : the URL of the article, PUBLISHER : the publisher of the article, CATEGORY : the category of the news item; one of: -- b : business -- t : science and technology -- e : entertainment -- m : health, STORY : alphanumeric ID of the news story that the article discusses, HOSTNAME : hostname where the article was posted, TIMESTAMP : approximate timestamp of the article's publication, given in Unix time (seconds since midnight on Jan 1, 1970),

Data Set link [1] (<http://archive.ics.uci.edu/ml/datasets/News+Aggregator>)

Document representation Techniques Used:

a. **bag-of-words (BOW):** is a simplifying representation used in natural language processing and information retrieval (IR). Also known as the vector space model. In this model, a text (such as a sentence or a document) is represented as the bag multiset of its words, disregarding grammar and even word order but keeping multiplicity. The bag-of-words model has also been used for computer vision. The bag-of-words model is commonly used in methods of document classification where the (frequency of) occurrence of each word is used as a feature for training a classifier. An early reference to "bag of words" in a linguistic context can be found in Zellig Harris's 1954 article on Distributional Structure.

b. Term frequency–inverse document frequency (TFIDF): Is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. It is often used as a weighting factor in searches of information retrieval, text mining, and user modeling. The tf-idf value increases proportionally to the number of times a word appears in the document and is offset by the frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general. Tf-idf is one of the most popular term-weighting schemes today; 83% of text-based recommended systems in digital libraries use tf-idf.

Classification Algorithms used :

Support Vector Machine (SVM): “Support Vector Machine” (SVM) is a supervised machine learning algorithm which can be used for both classification and regression challenges. However, it is mostly used in classification problems. In this algorithm, we plot each data item as a point in n-dimensional space (where n is number of features you have) with the value of each feature being the value of a particular coordinate. Then, we perform classification by finding the hyper-plane that differentiate the two classes very well .

Convolution neural network : Convolutional neural networks (CNNs) have dramatically improved the approaches to many active research problems. One of the key differentiators between CNNs and traditional machine learning approaches is the ability for CNNs to learn complex feature representations. We apply a CNN-based approach to categorization articles to (Business, science, health, entertainment). we building 3 layers all of them are dense layers, the first one take input with activation function “relu” with 512 neurons, the second one take 256 neurons with “sigmoid” activation function and finally the latest one is the output layer which take number of score to sum up the score of each class and take the highest one and used softmax activation function in last layer to make probability distribution and can easily know the win class

Logistic Regression (LG) :more general Framework , Logistic is a special mathematical function it uses and regression is Combines a weight vector with observations to create an answer More general cookbook for building conditional probability distributions . Naïve Bayes (later today) is a special case of logistic regression.

First split my data to train 80% and test 20%, this test determine even the training in right path or not it's called validation data. then this data have been tokenized and then converted to indices vector to can fit to train then we build 3 layers input hidden and output layer , each one have number of neurons , the first layer has 512 number of neurons , second layer contain 256 neuron and the last one have neuron to easily give score to each

Experiment #2(CNN):

Experiment #3(CNN):

Experiment #4 (CNN):

Experiment #5(KNN):

Experiment #6(KNN):

Experiment #7(KNN):

Experiment #8(KNN):

Experiment #9:

http://fci.helwan.edu.eg/wiki/%D9%85%D8%B9%D8%A7%D9%84%D8%AC%D8%A9_%D8%A7%D9%84%D9%84%D8%BA%D8%A7%D8%AA_%D8%A7%D9

As the experiment above but we changed the Max feature parameter to 10 with removing stop words.

As the experiment above but we use a Bag Of Words (BOW) representation of each document. We split the data, so that 20% of them remain for testing. And also we used Max feature as 1700 with removing stop words

In this experiment we use a Term Frequency (TF) representation of each document. And also a linear Support Vector Machine (SVM) classifier. We split the data, so that 20% of them remain for testing. As we changed the parameter use `idf` in `TfidfVectorizer` to `false`.

The same as the experiment #1 we changed the gamma parameter to 10 to fine tune the parameters to get the best results.

In this experiment we use word embedding as word and document representation technique with Recurrent neural network as Classification Model . First we splitting the data train and test (20% test) and preprocessing the data according to MAX_DOCUMENT_LENGTH we specified which is 10 and convert string label to and integer then Convert indexes of words into embeddings. This creates embeddings matrix of [n_words, EMBEDDING_SIZE=50] and then maps word indexes of the sequence into [batch_size, sequence_length, EMBEDDING_SIZE], after that Split into list of embedding per word, while removing doc length dim. Then Create an unrolled Recurrent Neural Networks to length of 10 as MAX_DOCUMENT_LENGTH and passes word_list as inputs for each unit. Given encoding of RNN, take encoding of last step and pass it as features for softmax classification over output classes.

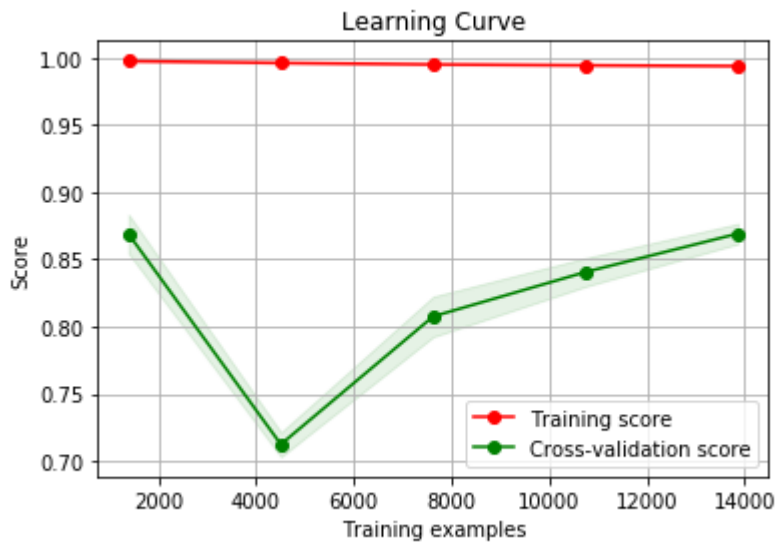
In this experiment we also use Recurrent neural network as Classification Model but with Bag of words (BOW) document representation . First we splitting the data train and test (20% test) and preprocessing the data according to MAX_DOCUMENT_LENGTH we specified which is 10 and convert string label to and integer then get feature column that represents sequences of integers based on n_words then Pass this to embedding_column to convert sequence categorical data into dense representation for input to sequence RNN ,then Builds input layer for given feature_columns , build last dense layer as classification layer then pass it as features for softmax classification over output classes .

In this experiment we use a TFIDF representation of each document. And also a Naïve Bayes (NB) classifier. We split the data, so that 20% of them remain for testing. And also we used Max feature as 170 with removing stop words .

As the experiment above but we changed the Max feature parameter to default with removing stop words.

As the experiment above but we changed the Max feature parameter to 9500 with removing stop words.

Activate Win



Experiment# 6 :

result is :

| Classification Report: | | | | | |
|------------------------|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| b | 0.90 | 0.91 | 0.90 | 1583 | |
| e | 0.93 | 0.96 | 0.94 | 1789 | |
| m | 0.91 | 0.88 | 0.89 | 551 | |
| t | 0.91 | 0.87 | 0.89 | 1277 | |
| avg / total | 0.91 | 0.91 | 0.91 | 5200 | |
| accuracy :0.941730769 | | | | | |

Graph of learning curve:



Experiment# 7 :

result is :

```
0.8813461538461539 : accuracy
```

The plot, titled "Learning Curve", shows the relationship between the number of training examples and the model's performance. The x-axis represents "Training examples" from 0 to 14,000, and the y-axis represents "Score" from 0.70 to 0.90. The red line (Training score) starts at approximately 0.93, drops to 0.78 at 4,500 examples, and then rises to 0.90 at 14,000 examples. The green line (Cross-validation score) starts at approximately 0.86, drops to 0.69 at 4,500 examples, and then rises to 0.85 at 14,000 examples. Both lines show a sharp initial drop followed by a gradual increase, with the training score consistently higher than the cross-validation score.

| Training examples | Training score | Cross-validation score |
|-------------------|----------------|------------------------|
| 1,500 | 0.93 | 0.86 |
| 4,500 | 0.78 | 0.69 |
| 7,500 | 0.86 | 0.79 |
| 10,500 | 0.88 | 0.82 |
| 14,000 | 0.90 | 0.85 |

result is :

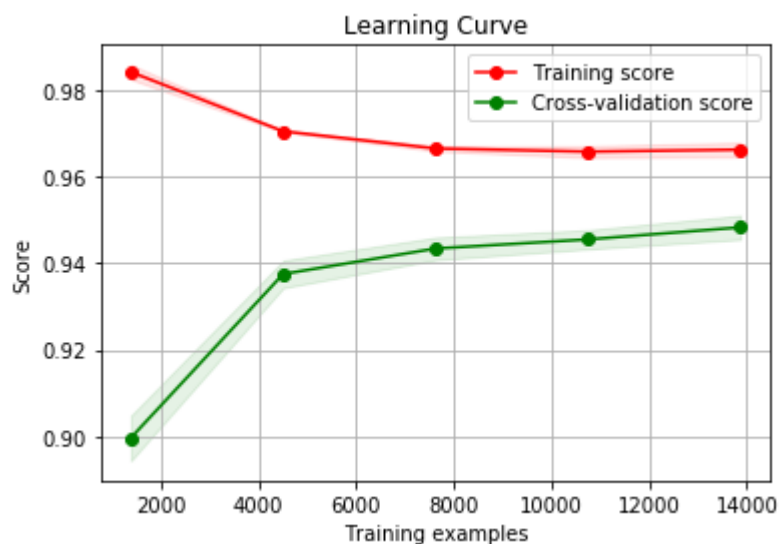
```
0.8948076923 : accuracy
```

Graph of learning curve:



Results:

The below figure is the learning graph of training and validation data .

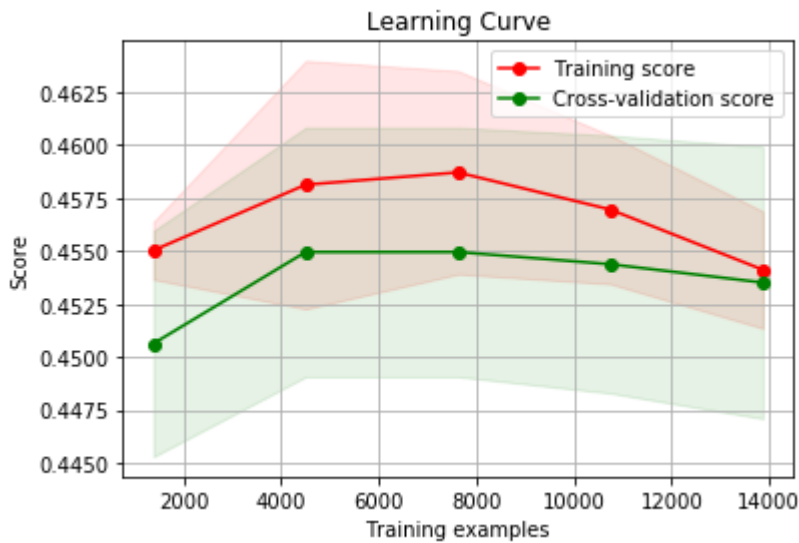


Results:

| Classification Report: | | | | |
|------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| b | 0.63 | 0.19 | 0.29 | 1583 |
| e | 0.39 | 0.98 | 0.56 | 1789 |
| m | 0.00 | 0.00 | 0.00 | 551 |
| t | 0.90 | 0.20 | 0.32 | 1277 |

```
avg / total      0.55      0.44      0.36      5200
Accuracy: 0.4428865178
```

The below figure is the learning graph of training and validation data .

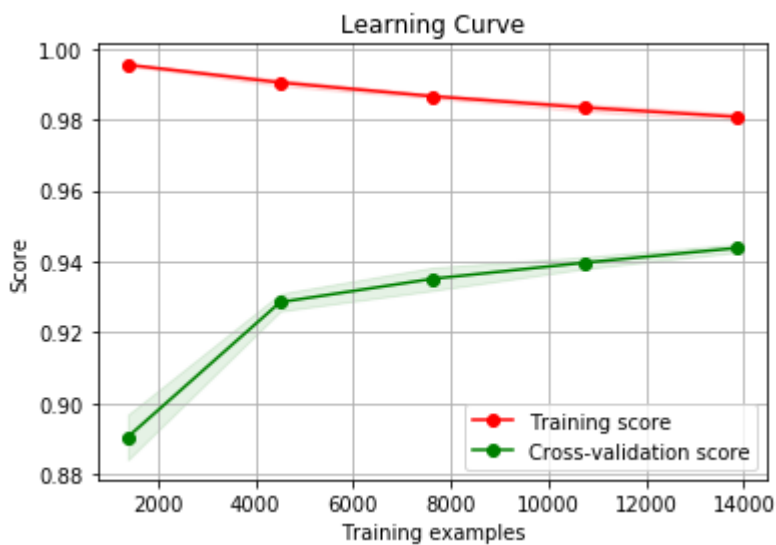


Experiment #11:

Results:

| Classification Report: | | | | |
|--------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| b | 0.91 | 0.96 | 0.93 | 1583 |
| e | 0.99 | 0.97 | 0.98 | 1789 |
| m | 0.96 | 0.91 | 0.93 | 551 |
| t | 0.94 | 0.93 | 0.93 | 1277 |
| avg / total | 0.95 | 0.95 | 0.95 | 5200 |
| Accuracy: 0.948269230769 | | | | |

The below figure is the learning graph of training and validation data .



Experiment #12:

Results:

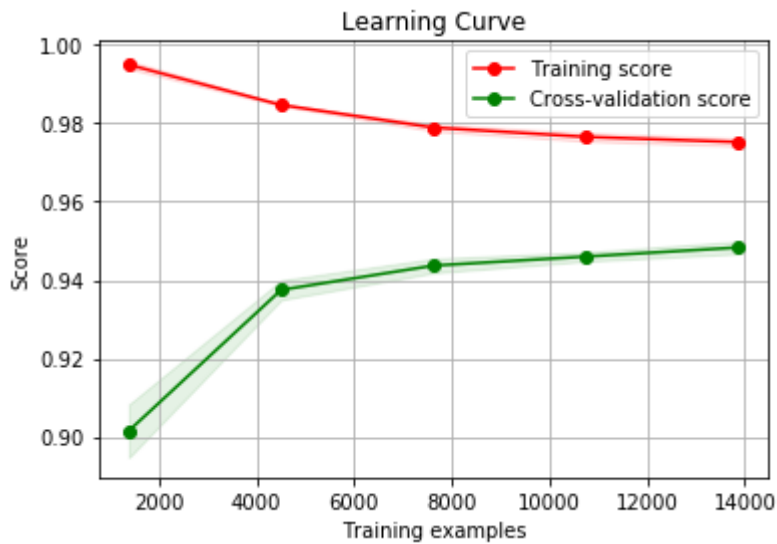
```
Classification Report:

```

| | precision | recall | f1-score | support |
|-------------|-----------|--------|----------|---------|
| b | 0.92 | 0.96 | 0.94 | 1583 |
| e | 0.98 | 0.97 | 0.98 | 1789 |
| m | 0.96 | 0.91 | 0.94 | 551 |
| t | 0.94 | 0.93 | 0.94 | 1277 |
| avg / total | 0.95 | 0.95 | 0.95 | 5200 |

```
Accuracy: 0.951153846154
```

The below figure is the learning graph of training and validation data .

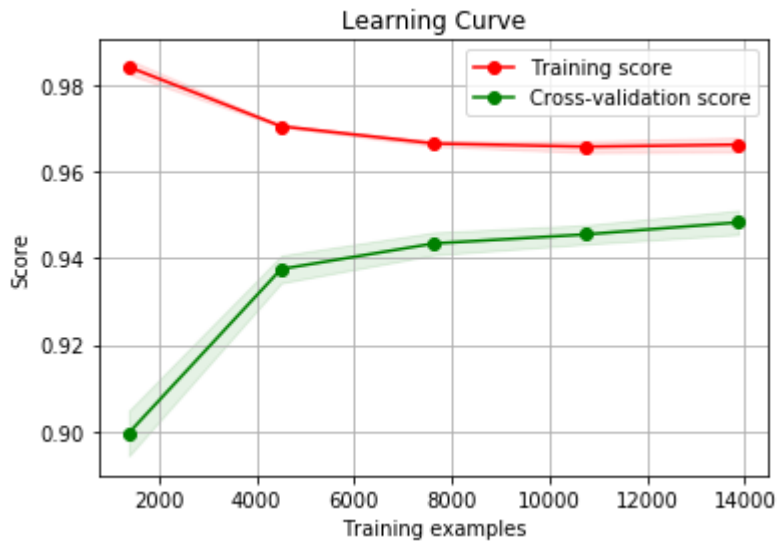


Experiment #13:

Results:

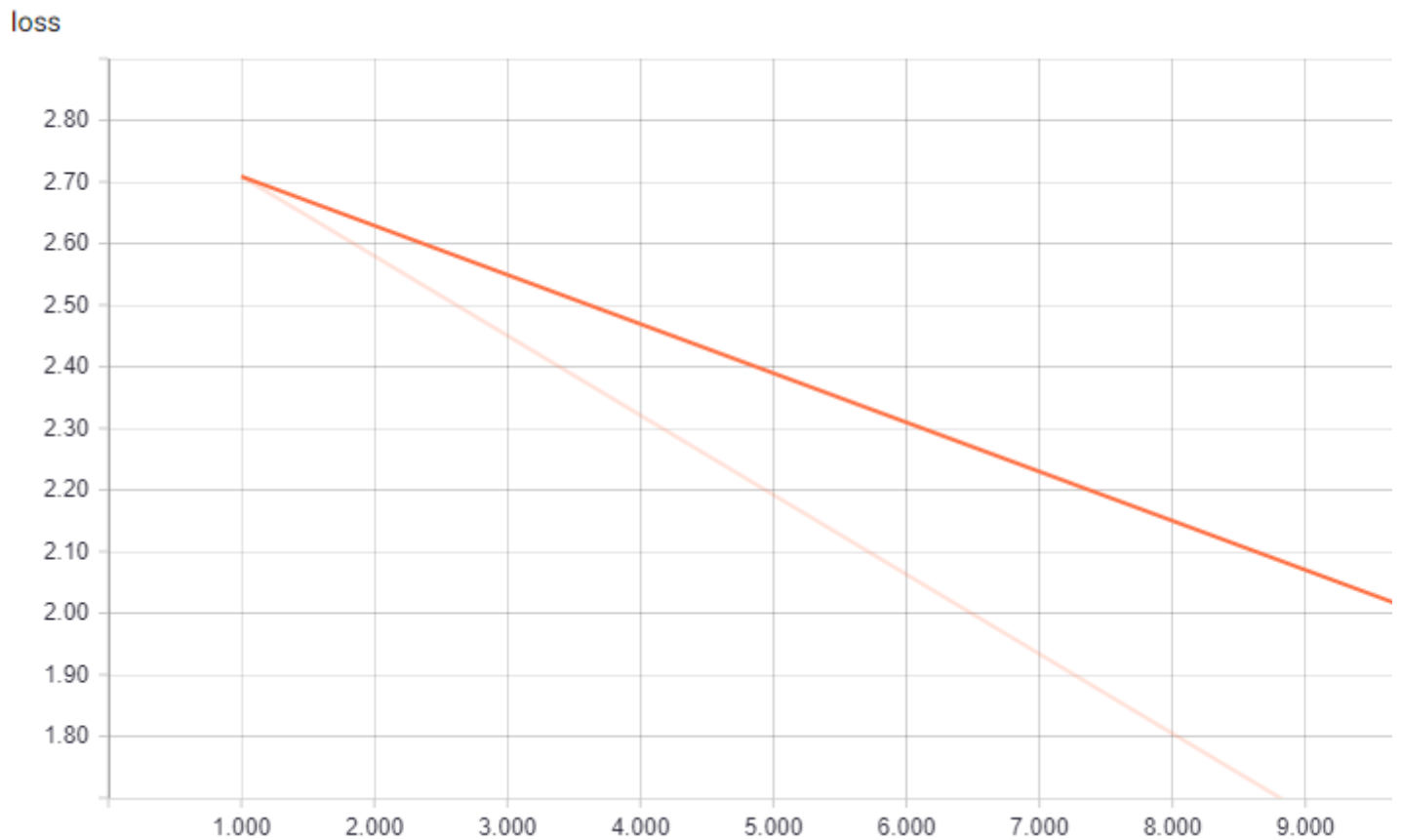
| Classification Report: | | | | |
|--------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| b | 0.90 | 0.97 | 0.93 | 1583 |
| e | 0.99 | 0.97 | 0.98 | 1789 |
| m | 0.98 | 0.90 | 0.94 | 551 |
| t | 0.95 | 0.92 | 0.93 | 1277 |
| avg / total | 0.95 | 0.95 | 0.95 | 5200 |
| Accuracy: 0.950576923077 | | | | |

The below figure is the learning graph of training and validation data .



Experiment #14:

Results: Loss graph:



```
Accuracy (tensorflow): 0.935000
```

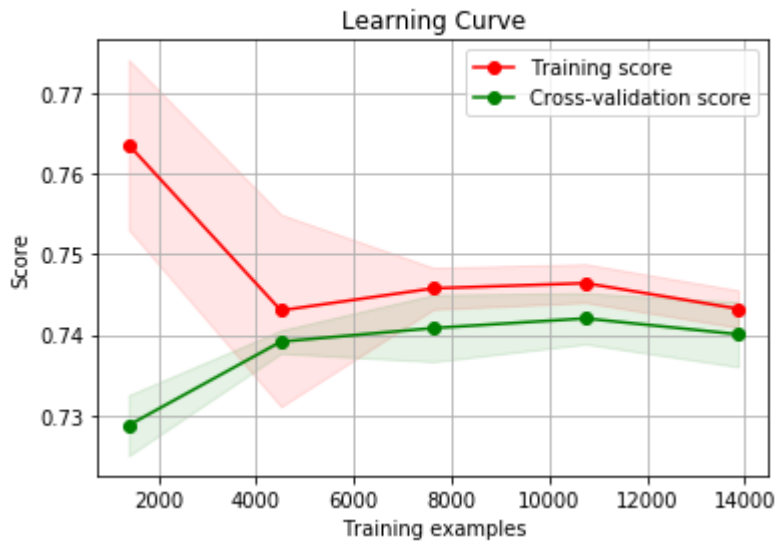
Experiment #15:

Results:

```
Accuracy (tensorflow): 0.943269
```

Experiment #16:

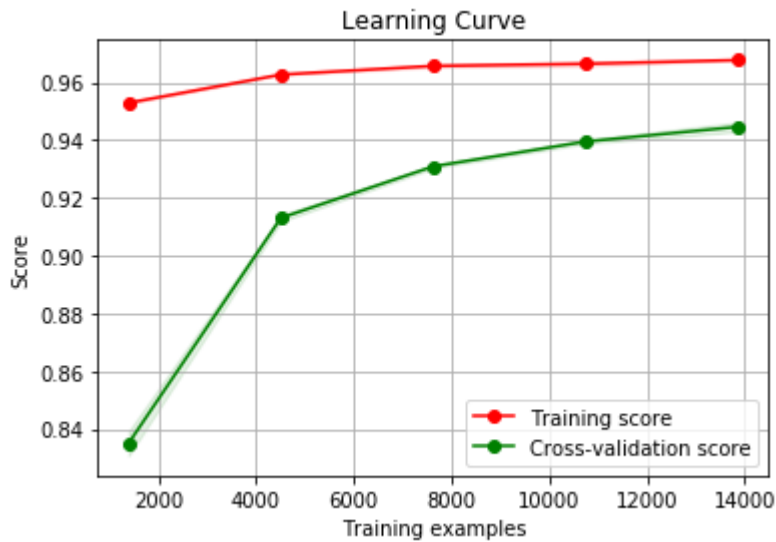
Results: Graph:



Accuracy : 0.732692307692

Experiment #17:

Results: Graph:



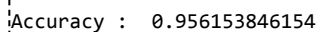
```
Accuracy : 0.955961538462
```

Experiment #18:

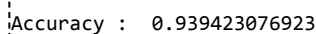
Results: Graph:



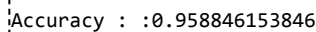
Results: Graph:



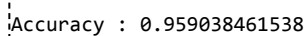
Results: Graph:



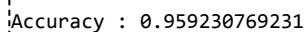
Results: Graph:



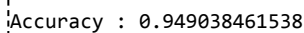
Results: Graph:



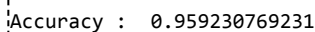
Results: Graph:



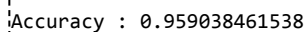
Results: Graph:



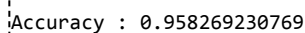
Results: Graph:



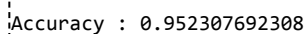
Results: Graph:



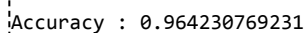
Results: Graph:



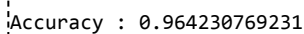
Results: Graph:



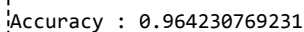
Results: Graph:



Results: Graph:



Results: Graph:



In Experiment#1, we used also adam optimizer and set out number of neurons to 512, 256, 64, 4 All of these give us a large percentage of accuracy in accuracy of training and validation

In experiment#4: when we decrease number of neurons the learning from data decrease but not with large scale

finally , Choosing the optimal K is almost impossible for a variety of problems , as the performance of a KNN classifier varies significantly when K is changed as well as the change of distance metric used and how they are uniformly distributed

In RNN model we experimented it with Word embedding and BOW document representation as we see the result above when experimented the RNN with bag of words and word embedding it seems that word embedding result best due to word embedding is an improvement over more the traditional bag-of-word model encoding schemes where large sparse vectors were used to represent each word or to score each word within a vector to represent an entire vocabulary is replaced with dense vectors where a vector represents the projection of the word into a continuous vector space.

http://fci.helwan.edu.eg/wiki/%D9%85%D8%B9%D8%A7%D9%84%D8%AC%D8%A9_%D8%A7%D9%84%D9%84%D8%BA%D8%A7%D8%AA_%D8%A7%D9

appropriate techniques, in this case, the best result of this experiment when we make max_feature = 9500. From experiment 20 to 23 we use NB with BOW and exchange the max_Feature as a four experiment above the best accuracy appear when make the Max_Feature = 11500.

In LG algorithm we use TF-IDF or BOW and calculate the accuracy according to validation data. Form experiment 24 to 27 we use LG with TF-IDF and exchange the max_Feature as a four experiment above the best accuracy appear when make the Max_Feature is default value. form experiment 28 to 31 we use LG with BOW and exchange the max_Feature as a four experiment above the best accuracy appear when make the Max_Feature =11500.

Conclusion

In this project we worked on News Classification problem we experimented many algorithms K-nearest neighbor , Logistic regression , Naïve bayis , Support vector machine and Neural network models like CNN and RNN, we used different document representation like Bag of words , TF/IDF and Word embedding . All the training done on the News-aggregator data set which is labeled data set, After trying different models it turns out the accuracy ranges from 34% to 99% depends on hyper parameters in each model. we find CNN model in experiment 1 is better than other due to its accuracy that we recommended.

مجلوبة من "http://fci.helwan.edu.eg/w/index.php?title=7/مشاريع/2018/معالجة_اللغات_الطبيعية&oldid=611"

■ آخر تعديل لهذه الصفحة كان يوم ١ مايو ٢٠١٨ الساعة ١٠:٤٦.