

## SENTIMENT ANALYSIS ON SOCIAL MEDIA TWEETS

A.P.Uthayashankar

unstoppable2397@gmail.com

UG Scholar, Amrita Vishwa Vidyapeetham University, Coimbatore, Tamilnadu,India

### Abstract

Due to the large amount of data present in the internet, companies can use these data to their advantage to increase their profit. The way that the companies can use this data to their advantage is through Sentiment Analysis, which in simple terms is basically called opinion mining, i.e. collect the opinion of people on a particular subject. It classifies the tweet into positive, negative and neutral. These opinions collected from people can be very helpful for big companies to make big and important decisions correctly which will help them to develop further. In this work, a gradient descent based RMSprop optimization algorithm for the classification of tweets as either positive, negative or neutral. RMSPROP stands for Root Mean Square prop which is an algorithm that is used to speed up the gradient descent of a CNN. CNN stands for Convolutional Neural Network, which is a main category in deep learning that is used for classifications, detections, and recognition of objects. CNN takes an input, processes it by passing the input through various layers of the CNN where each layer performs a unique function and finally it classifies the input in a certain category as the output of the CNN. Boosting the gradient descent of a CNN helps in achieving higher accuracy of the result since gradient descent helps to improve deep learning and CNN models by reducing the cost function. In this work predicting the type of opinion by RMSprop based CNN network is studied. All the tweets undergoes pre-processing and feature extraction process for efficient classification of the tweets. The performance evaluation of the proposed method is tested in MATLAB Tool, which achieves the most efficient classification accuracy.

**Keywords:** Tweets classification, Foreign Languages, CNN, gradient descent, RMSprop, Accuracy.

### 1. Introduction:

In the present day, millions and trillions of people worldwide have social media accounts and they express their opinions on all types of subject matter. They always have something or the other to express their opinion about something. Therefore social media is the best source of data that companies can rely on to make important decisions for their future development. They can also do

sentiment analysis on the data present in web pages like amazon and flipkart where people give product reviews about a lot of products. The respective company can do sentiment analysis on their product's product review to mine for opinions on their products which can help the growth of their business. Sentiment analysis is not only used in e-commerce but also to make important political decisions. For example a politician can make a post in the social media about a particular topic and wait for the people to respond. Now these responses can be collected and sentiment analysis can be done on them, depending on these opinions, the politicians can make their next decision accordingly.

An analysis about the sentiment expression based on the machine learning algorithms. Emotions are used in social media to express their opinion on a subject. Based on the polarity of the words which express different emotions [1].

A binary classification technique for sentiment analysis for twitter data by utilizing distant supervision, in which their training data consisted of tweets with emojis which filled as noisy labels [2]. Naive Bayes (NB) and Support Vector Machines (SVM) are mainly utilized for classification. The performance results demonstrates that SVM achieves the most efficient results when compared to NB and unigram models obtain more efficient results when compared to bigram models for feature extraction.

A Twitter API to collect twitter data and classify them as either positive or negative tweets [3]. The algorithm is applied to the training data set which filters the opinions based on the tweets' contents. Unigram NB is used to simplify the data by eliminating features based on mutual information. Feature extraction by chi square and the algorithm is applied to classify them as positive or negative tweet.

An efficient method for classification of opinion by integrating certain features [4]. The Part-of-speech information and Word-relations are used as feature sets for classifying the positive and negative tweet based on certain classifiers NB, Maximum Entropy (ME) and Support Vector Machines. The performance results show that method under study achieves higher accuracy.

A supervised sentiment classification system to separate the hash tags and smileys. The 50 Twitter tags and 15 smileys are used as sentiment labels [5]. The different kinds of sentiment classification are used to recognize untagged sentences. The performance results of the quality of the sentiment identification are assessed by human judges. The expected results indicates that the efficient classification of the smileys is based on cross validation.

A Naive Bayes for sentiment analysis for the recognition of the English tweets [6]. This classifier is used to classify the positive, negative and neutral tweets. The model performs in basic rule for polarity words for analysing tweets. The data set is pre-processed to remove the noise and background illumination and then it is passed on to the Naive Bayes classifier for classifying the tweets. The experimental results show that the classifier accomplishes the effective performance of 63% f-score for tweets polarity classification.

A machine learning technique based on computer-assisted techniques for effective analysis of sentiment in [7]. The pre-processing stage is executed to remove the punctuation, tags and transform them into structured form. The lexicon techniques are employed to accomplish numerical score and the feature selection utilized for sentiment analysis. The SVM and RBF kernel is applied for the classification, which achieves the efficient classification accuracy.

A model for tax comments analysis using text mining is proposed in [8]. The input data set are collected from the Facebook and twitter for processing tax comments. Initially pre-processing is applied to the transformation of structured format, the mining of text involves the phases of tax comments, the feature selection is implemented to identify the appropriate features, and the SVM for analysing the tax comments does classification.

In this, a deep learning based on the gradient descent algorithm for the classification of foreign language tweets is proposed.

The paper is organised by explaining the existing techniques for sentimental classification in section 2. The section 3 describes about the proposed method. Section 4 describes about the implementation procedure and discussion based on the results. The paper is concluded with conclusion in section 5.

## **2. Related Works:**

The survey of sentimental classification of tweets is discussed in [9]. The paper described about the different twitter datasets, different

feature extraction methods and classification methods base line algorithm and naïve Bayes classification.

A part of speech based classification using the Naïve Bayes classifier is performed for the twitter datasets [10]. The datasets undergoes url removal and part of speech extraction process. The classifier is trained based on the POS and its saliency and entropy value for tweet classification.

A tree based classification is used for the classification of tweets in [11]. Here, the tweets are converted into sentence format using emotion dictionary and acronym dictionary. Then the data is processed and part of speech is extracted. The part of speech is grouped in the form of tree with root nodes as positive, negative and neutral.

A different classification methods namely SVM, naïve Bayes, Random forest, SMO and filtered classification is used for opinion mining in [12]. The pre-processing steps for all the classifier is same. The tweets are classified as positive and negative using different classifiers. The performance is evaluated is based on the accuracy, precision and recall. Among different classifiers, the Bayesian logistic Regression classifier produced the best result with the accuracy of 75 %.

The classification of Hindi tweets using SVM and decision classifier is discussed in [13]. The pre-processing step involves user name removal and hyper link removal. The feature extraction is based on the TIFIDF score of unigram vector for both the dictionary and the testing vector. The classifier is trained with the scores of dictionary and tested with the scores of testing tweets using WEKA. SVM classifier shows the best results as compared to the decision classifier J48.

The classification of Italian tweets for Amazon EBook review is discussed in [14]. Here, the lexicon is built for the tweets by converting the Italian tweets to English language. Then the feature extraction process is performed the pre-processed tweets. The feature extracted tweets are classified using AOL based on the scores.

The Classification of Arabic tweets using external source called SentiwordNet using machine learning approach is discussed in [15]. In this, the bag of words is extracted for the Arabic tweets. The scores are generated for the corresponding English version and Arabic version of Bag of words. Then, the scores used

for classifying the tweets using machine learning algorithm.

### 3. Proposed Method

The Sentiment Analysis system for social media tweets that may be helpful to analyse tweets that are highly unstructured in their meaning using CNN model and optimizing it with RMSprop optimization algorithm and identify the opinion they express as either positive, neutral or negative. The sentiment analysis is also done on three different languages namely Spanish, French and Dutch.

Sentiment Analysis in twitter is complicated because of minimum length format of less than or equal to 100 words. Generally people share their opinions in the form of unstructured representation which includes blogs, forums etc. In this model, the analysis is made to differentiate the twitter tweets by two steps, first one is to remove the unstructured blogs and tags and in the second classify the tweets by learning algorithms. Here are some of the processes that takes place in the sentiment analysis model.

#### 3.1 Data Acquisition:

The tweets (data set) are collected from the publicly available data sets of twitter link. Data in the form of raw tweets is obtained from the internet which provides a sufficient amount of data to test and train the CNN model.

#### 3.2 Pre-processing:

The pre-processing is mainly used for cleaning, background homogenization and reduce the noise in the input image to enhance the image quality. During data acquisition and transmission the image is corrupted by noise, so it is necessary to reduce noises by preserving the significant image features.

The main process in pre-processing step is the removal of emoji's and URL. A tweet contains an opinion about a subject which are expressed in different ways by different people. The raw data having noise is highly affected by inconsistency and redundancy. The main process of pre-processing steps is that the URL and tags are removed, corrects the spelling and repeated characters, punctuations, symbols and numbers are removed.

#### 3.3 Feature Extraction:

This goal is to extract the known feature, which contains some clue about the opinion of the tweet. In this process the known feature is extracted, which contains some details about the opinion of the tweet. The various feature extraction techniques are available to gather the relevant features from text, which can be extracted in two steps. In the first step, twitter specific features are removed to create normal text. In the second step the feature extraction is implemented to get the appropriate features. The pre-processed data set consists of unique properties that are forwarded to the feature extraction method to extract the aspects of the processed data to classify as either positive, negative or neutral.

The twitter tweets data sets are collected from various resources and are pre-processed to extract the features that will help determining the opinion of each tweet. The steps involved in feature extraction are tokenization, stop word removal, stemming and N-gram.

##### i. Tokenization:

Tokenization is the process, which is used to break down a sentence into words, phrases, symbols or other meaningful tokens by extracting and eliminating the punctuation marks from the given sentence.

##### ii. Stop Word Removal:

Stop words are the commonly used words which appear in any language which are not useful for determining the opinion of the sentence and hence they are removed for focusing on the more important words which will help in determining the opinion of the sentence or tweet. For example in English, the examples of stop words are the, is, and, about, above, of, at, etc..

##### iii. Stemming:

Stemming is the process of transforming a word into a more simple form by removing the suffixes that a word ends with. For example the suffixes are ed, ion, inos, ing, le etc.. The words in this format of a sentence can help to determine the opinion of the sentence more easily.

For example automate, automatic and automation are all reduced to automate

##### 1. Porter Stemming:

Porter Stemming is the process of removing suffixes from English words using some defined functions. This process is mainly used for retrieving information. It is represented by a

group of words or certain terms. This algorithm is generated by certain assumptions without the stem dictionary to improve the English sentence.

The goal is to remove the suffixes from two words R1 and R2 to generate the single stem S. In some cases if there is no difference between the two statements, then it is represented as “a document about R1” and a document about R2”. For an instance R1= CONNECTION and R2 = CONNECTIONS are defined into single stem since R1 and R2 are not entirely different and have the same meaning.

But if R1 = RELATE and R2 = RELATIVITY are defined as different and stemmed into two separate because R2 is related with physics and R1 is not. Both the cases are different.

The rules to remove a suffix is given as :

(Condition) R1 to R2

It defines that the word ends with suffix R1 and the stem before R1 satisfies the condition and R1 replaced by R2.

For example the word ends with EMENT  
(m>1) EMENT

R1 is EMENT and R1 is null. This condition map the REPLACEMENT to REPLAC

This algorithm is implemented not only to eliminate the suffix even when the stem is too short and the length of the stem is measured to be m but also this method is used for efficiently removing the suffix.

### 3.4 N-GRAM:

N-gram defines the combinations of neighbouring words or length of letters n in the given tweet. An n-gram represents the group of n words or characters (indicated as grams denotes grammar) which follows one another. N-grams can be used to predict the next word given the previous N -1 word. N-grams are widely used in data mining and word processing tasks.

For an instance If “I work in India “ denotes n=4 (number of words from the sentence which is used to create an index of how often words follows one another.)

N-gram can be represented as:

$N\text{-grams}_k = X-(N-1)$

The X represents number of words in given sentence of a tweet and N denotes number of grams for K sentence. N-grams are widely used for various word processing tasks. The language model is generated by using n-grams to develop unigram, bigram and trigram models.

The fundamental operation of N-gram that captured the language representation from arithmetic view indicates that the letter or word is given as the input structure. If N-gram is larger i.e. n is higher then it denotes more number of text. The length of the word can be represented based on the application. If n-grams are minimum, then it fails to detain the differences. Suppose n-grams are long, then it fails to detain the particular cases. The figure 1 shows an example for N-gram.

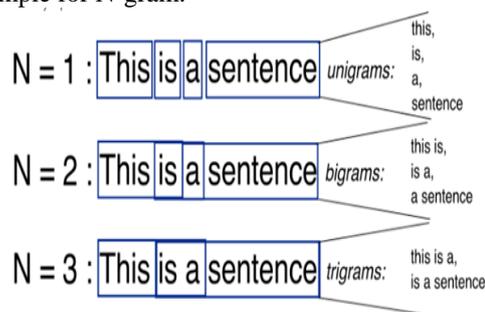


Fig 1. N-Gram Classification

### 3.5 Classification:

The extracted feature values from the N-gram are forwarded to the classifier to distinguish the positive, negative and neutral tweets. The learning algorithms are mainly used to solve the classification problems. The raw data sets undergoes training by classifier to classify the positive, negative and neutral tweets. Convolutional Neural Network is used here as the classifier to classify the positive, negative and neutral tweets [16]. The CNN contains neurons in three dimensions, which are spatial dimensionality of the input layer and depth.

The CNN consists of three layers, they are convolutional layer, pooling layer and fully connected layer. CNN is formed by stacking these layers

#### 3.5.1 Pooling Layer:

It minimizes the dimensionality of the representation and also minimizes the amount of parameters needed in the model and the computational complexity of the model.

#### 3.5.2 Fully Connected Layer:

This layer includes certain neurons that are directly connected to the neurons in the neighboring layers, without connecting to any other layers.

The main function of CNN

- The input layer maintains the input image’s pixel value.

- The convolutional layer identifies the neurons output which are linked to the input by weights. The previous layer generates the activation functions, which are sigmoid to the output.
- The pooling layer operates by reducing the sampling with the spatial dimensionality of the input layer as it reduces the number of parameters.
- The main function of the fully connected layer is to generate the activation function which will be used for classifying the tweets.

When the input data is transferred to the convolutional layer, then the each layer convolves the filter around the spatial dimensionality of the input to generate the 2D activation map. CNN comprises of neurons that adapt themselves by learning. Every neuron will receive input and operated based on scalar product of non-linear function. For CNN operation, the input and the output layer are expressing a single activation function (weight). The final layer applies the loss function linked with the activation weight.

The working of RMSprop optimization based CNN for sentimental analysis is given below:

1. The data set for twitter tweets is taken from Stanford University.
2. The data set taken in this work consists of three foreign languages and English language.
3. These tweets are then classified as training and testing tweets using cross-validation using holdout approach.
4. The training tweets undergoes the pre-processing steps like URL removal, stop word removal and stemming process.
5. The pre-processed tweets are then converted into tokenized documents to store the important words which give details of the opinion from the tweets.
6. These words are used for training the network.
7. In this only one set of the convolution neural network layers are used for the training.
8. The CNN layers consists of the input layer, the convolution layer, batch normalization layer, RELu layer, max-pooling layer, fully connected layer, soft-max layer and classification layer.
9. The layer is trained through the RMSprop optimization, which improves the learning

- rate of the model and it is adaptive in nature.
10. The trained CNN is saved for classifying the tweets.
11. The testing data is undergoes the same preprocessing and feature extraction process.
12. The extracted features are given as the input to trained network to classify the tweets.

The evaluation of the proposed method is based on the following metrics in section 2.5

### 3.5 Performance Analysis:

The performance of our proposed CNN method is simulated in MATLAB user interface under the windows environment. The evaluation of the proposed method is measured by certain performance metrics such as accuracy, sensitivity, specificity and F-measure, which is obtained by using confusion matrix.

#### 3.5.1. Accuracy:

Accuracy measures the exact classification of tweets comments

$$Accuracy = (TP + TN) / (TP + TN + FN + FP) \quad 1$$

#### 3.5.2 Specificity:

Specificity measures the negative classification of the tweet when the condition is actually not present. It is recognize as false-positive rate.

$$Specificity = TN / (TN + FP) \quad 2$$

#### 3.5.3 Sensitivity:

Sensitivity measures the positive classification of tweet comments when the condition is actually present. It is represented as false-negative rate.

$$Sensitivity = TP / (TP + FN) \quad 3$$

Where

- True Positive (TP) means positive result of the tweet classification
- True Negative (TN) means negative result of the tweet classification

- False Positive (FP) means the positive result of the negative tweets classification
- False Negative (FN) means the negative result of the positive tweets classification

labeled data sets are then analysed by pre-processing and feature extraction techniques. Pre-processing is applied to raw data sets and deep learning networks are applied to train the data set with the feature vectors from the large data set to classify the twitter data sets as either positive, negative or neutral.

#### 4 Results and discussion

The proposed method RMS-CNN is used for the classification of tweets based on the sentiments for foreign languages. The whole process is tested using the MATLAB software with R2018a version under windows 10 environment.

The data sets are collected from the publicly available resources from cnlp\_language processing. The French and the Spanish data sets is formed from the labeled English data sets by translating it. A sample of labelled data are shown in Fig 2. The

Var1	Var2
"rt brax ying : ma fille morganwhitlfly fait incroyable dans la pièce :"	3
"rt paulabroadwell : conseil de réseautage pour aujourd'hui : payez-le à l'avance !"	3
"cette an née john lewis annonce est tellement merde par rapport aux dernières an nées"	3
"jodiechadderton ce qui est mal , je peux voir la similitude! si quelque chose je pense que c'est une amélioration jo"	3
"pommes , pas de caféine , sont plus efficaces à vous réveiller le matin . intéressant didyouknow"	3
"rt kgejerta : les requins ont une semaine qui lui est dédiée. il est l'homme le plus intéressant dans le monde . mrpi2012 teamtim"	3
"länzicrellin tom je ve vous a obtenu 23 cadeaux pour votre anniversaire ha ami gentil"	3
"juste reconnaissant que je m ici"	3

Fig 2 Labelled data of the tweets

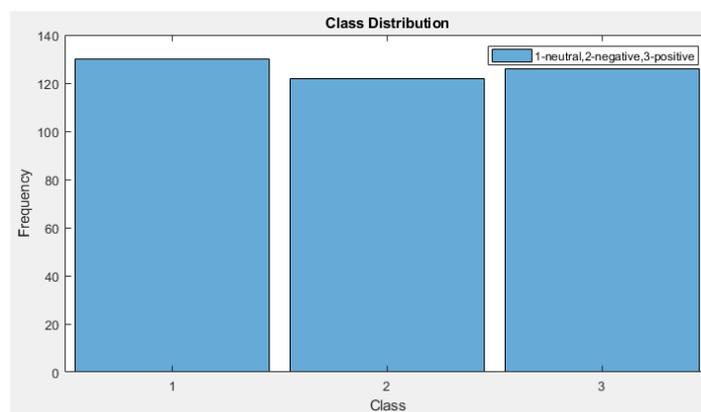


Fig 3. Distribution of tweets as positive(3), neutral(1) and negative(2)

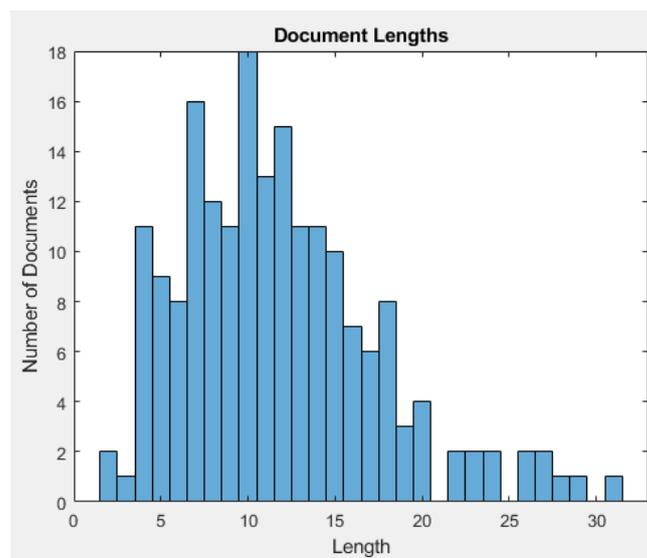


Fig 4. Distribution of tokenized documents in the set of tweets

```

Command Window
New to MATLAB? See resources for Getting Started.
Columns 99 through 100

(100x75 single) (100x75 single)

Training on single GPU.
Initializing image normalization.
=====
| Epoch | Iteration | Time Elapsed | Mini-batch | Mini-batch | Base Learning |
|       |          | (hh:mm:ss)  | Accuracy  | Loss       | Rate          |
|-----|-----|-----|-----|-----|-----|
| 1     | 1       | 00:00:00    | 36.72%    | 1.0986     | 0.0010       |
| 50    | 50      | 00:00:02    | 58.59%    | 0.7758     | 1.0000e-07   |
| 75    | 75      | 00:00:03    | 58.59%    | 0.7758     | 1.0000e-10   |
|-----|-----|-----|-----|-----|-----|

convnet =
SeriesNetwork with properties:
    Layers: [20x1 nnet.cnn.layer.Layer]

YPred =
    
```

Fig 5. CNN properties of sentiment analysis

To extract the important features from the tweets, the stemmed words are transformed into tokenized documents. Length of tokenized document refers to number of words a sentence/tweet(document) has been split into for feature extraction purpose. The distribution of the tokenized documents are shown in Fig 4

This tokenized document is transformed into a sequence for the training of deep learning CNN. The CNN properties are shown in the above figure on page

The data set is classified into two types, namely training and testing data. The training data is trained through the CNN model with the above properties shown in Fig 5 and the network is

preserved for testing purpose. The evaluation of the proposed model is based on the mentioned metrics and the accuracy of the proposed model is compared with the existing method for single language.

The proposed RMSPROP optimization based classifier produced the best results because the weights for outputs are assigned in adaptive manner and its learning rate is incremented compared to the other optimization models. It also implemented all type of layers in the CNN model to improve the feature for training the network and with only 1 set of CNN the sentimental analysis of tweets performed well. The result that is obtained from the CNN model of sentiment analysis is the overall opinion of the sentence or tweet. There are three classifications of opinions, they are positive opinion, negative opinion and neutral opinion. Therefore the

CNN model will give an output of either 1, 2 or 3. Here 1 denotes neutral opinion, 2 denotes negative opinion and 3 denotes positive opinion.

The sample output results for each of the opinion type i.e. positive opinion denoted by 3, negative opinion denoted by 2 and neutral opinion denoted my 1.

```

77 - textDataTest='nou wat m n ma net zij moet ik jeemiselee morgen echt ff vertellen face2face hahahaha zo niet normaal ppp';
Command Window
New to MATLAB? See resources for Getting Started.

Columns 85 through 91
{100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 si

Columns 92 through 98
{100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75 si

Columns 99 through 100
{100x75 single} {100x75 single}

Training on single GPU.
Initializing image normalization.
=====
| Epoch | Iteration | Time Elapsed | Mini-batch | Mini-batch | Base Learning |
|       |          | (hh:mm:ss)  | Accuracy  | Loss       | Rate          |
|-----|-----|-----|-----|-----|-----|
| 1     | 1       | 00:00:00    | 28.91%   | 1.0986    | 0.0010       |
| 50    | 50      | 00:00:02    | 64.84%   | 0.8052    | 1.0000e-07   |
| 75    | 75      | 00:00:03    | 64.84%   | 0.8052    | 1.0000e-10   |
|-----|-----|-----|-----|-----|

convnet =
SeriesNetwork with properties:
    Layers: [20x1 nnet.cnn.layer.Layer]

YPred =
categorical
3
    
```

Fig 6 Positive opinion of a tweet

Here the input sentence is shown in the topmost part of the above image. The variable storing the input sentence is “textDataTest”. We can see the output in the bottom most part of the image under a sub heading called “categorical”. Under this sub heading you can see the number “3” which

denotes that the sentence contained the variable “textDataTest” expresses a positive opinion. Similarly the output screenshots are given for negative and neutral opinions are shown in Fig 7 and Fig 8 respectively

```

77 - textDataTest='makinnlovee suena como si estamos en el mismo barco : crappystart';
Command Window
New to MATLAB? See resources for Getting Started.
Columns 85 through 91
{100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75
Columns 92 through 98
{100x75 single} {100x75 single} {100x75 single} {100x75 single} {100x75
Columns 99 through 100
{100x75 single} {100x75 single}

Training on single GPU.
Initializing image normalization.
=====
| Epoch | Iteration | Time Elapsed | Mini-batch | Mini-batch | Base Learning
| | | (hh:mm:ss) | Accuracy | Loss | Rate
=====
| 1 | 1 | 00:00:00 | 31.25% | 1.0986 | 0.0010
| 50 | 50 | 00:00:02 | 57.81% | 0.8118 | 1.0000e-07
| 75 | 75 | 00:00:03 | 57.81% | 0.8118 | 1.0000e-10
=====

convnet =
SeriesNetwork with properties:
Layers: [20x1 nnet.cnn.layer.Layer]

YPred =
categorical

2
    
```

Fig 7. Negative opinion of a tweet

Under this sub heading called “categorical” that the sentence contained the variable in Fig 7, you can see the number “2” which denotes “textDataTest” expresses a negative opinion

```

77 - textDataTest='k ruim morgen me kamer wel op';
Command Window
New to MATLAB? See resources for Getting Started.
Columns 85 through 91
{100x75 single} {100x75 single} {100x75 sing
Columns 92 through 98
{100x75 single} {100x75 single} {100x75 sing
Columns 99 through 100
{100x75 single} {100x75 single}

Training on single GPU.
Initializing image normalization.
=====
| Epoch | Iteration | Time Elapsed | Mini-batch
| | | (hh:mm:ss) | Accuracy
=====
| 1 | 1 | 00:00:00 | 37.50
| 50 | 50 | 00:00:02 | 61.72
| 75 | 75 | 00:00:02 | 61.72
=====

convnet =
SeriesNetwork with properties:
Layers: [20x1 nnet.cnn.layer.Layer]

YPred =
categorical

1
    
```

Fig 8. Neutral opinion of a tweet

Under this sub heading called “categorical” in Fig 8, you can see the number “1” which denotes As we saw in the above images we can test any sentence by storing the sentence in the variable “textDataTest” and perform sentiment analysis on it to get the output either as 1(neutral opinion), 2(negative opinion) and 3(positive opinion). The sentence can be in any one of the four languages namely Spanish, French, Dutch and English.

The CNN model achieves a highest accuracy of 81.75% which means that out of all the tweets that need to be correctly classified, the CNN model has classified 81.75% of all the tweets correctly. It achieves the sensitivity of 58.71% which means that out of all the tweets that need to be correctly classified as true, the CNN has classified 58.71% of those tweets correctly. It achieves the specificity of 80.19% which means that out of all the tweets that need to be correctly classified as false, the CNN has classified 80.19% of those tweets correctly. Therefore from this testing we have successfully studied and confirmed that the CNN model using RMSPROP algorithm using gradient descent is an efficient model to perform the most effective sentiment analysis on social media tweets since even after testing with multiple number of data sets.

#### 4. Conclusion

In this project we have successfully studied a sentiment analysis model based on the CNN method using RMSPROP algorithm using gradient descent identify the type of opinion of a sentence. The sentence may be in any one of the four languages namely Spanish, French, Dutch and English. It successfully classifies any sentence into either positive opinion, negative opinion or neutral opinion, it displays the output as 3, 2 or 1 respectively. This CNN model is trained to classify social media tweets in either Spanish, French, Dutch or English.

To summarize the overall process, first the social media tweets data sets are pre-processed to remove hash tags and punctuations to enhance the tweet observations. Then feature extraction process is employed by stemming and N-gram. The RMSprop based Convolutional neural network is employed to classify the tweets as either positive opinion, negative opinion or neutral opinion.

After testing multiple data sets with the existing method of CNN using RMSprop algorithm, we see that it achieves a maximum accuracy of 81.75%, a maximum sensitivity of

that the sentence contained the variable “textDataTest” expresses a neutral opinion. 58.71% and a maximum specificity of 80.19%. Hence we can conclude that CNN using RMSprop algorithm is an efficient model to perform sentiment analysis.

#### References:

1. Wang, H. and Castanon, J.A., 2015. Sentiment expression via emoticons on social media. arXiv preprint arXiv:1511.02556.
2. Liang, P.W. And Dai, B.R., 2013, June. Opinion Mining On Social Media Data. In Mobile Data Management (Mdm), 2013 Ieee 14th International Conference On (Vol. 2, Pp. 91-96). Ieee.
3. Go, A., Bhayani, R. And Huang, L., 2009. Twitter Sentiment Classification Using Distant Supervision. Cs224n Project Report, Stanford, 1(12).
4. Xia, R., Zong, C. And Li, S., 2011. Ensemble Of Feature Sets And Classification Algorithms For Sentiment Classification. Information Sciences, 181(6), Pp.1138-1152.
5. Davidov, D., Tsur, O. And Rappoport, A., 2010, August. Enhanced Sentiment Learning Using Twitter Hashtags And Smileys. In Proceedings Of The 23rd International Conference On Computational Linguistics: Posters (Pp. 241-249). Association For Computational Linguistics.
6. Citius: A Naive-Bayes Strategy For Sentiment Analysis On English Tweets. In Proceedings Of The 8th International Workshop On Semantic Evaluation (Semeval 2014) (Pp. 171-175).
7. Jadav, B.M. And Vaghela, V.B., 2016. Sentiment Analysis Using Support Vector Machine Based On Feature Selection And Semantic Analysis. International Journal Of Computer Applications, 146(13).
8. Utami, E. And Luthfi, E.T., Text Mining Based On Tax Comments As Big Data Analysis Using Svm And Feature Selection.
9. Kharde, V., & Sonawane, P. (2016). Sentiment analysis of twitter data: a survey

- of techniques. arXiv preprint arXiv:1601.06971.
10. Pak, A., & Paroubek, P. (2010, May). Twitter as a corpus for sentiment analysis and opinion mining. In *LREc* (Vol. 10, No. 2010, pp. 1320-1326).
  11. Agarwal, A., Xie, B., Vovsha, I., Rambow, O., & Passonneau, R. J. (2011, June). Sentiment analysis of twitter data. In *Proceedings of the workshop on language in social media (LSM 2011)* (pp. 30-38).
  12. Gokulakrishnan, B., Priyanthan, P., Ragavan, T., Prasath, N., & Perera, A. (2012, December). Opinion mining and sentiment analysis on a twitter data stream. In *International Conference on Advances in ICT for Emerging Regions (ICTer2012)* (pp. 182-188). IEEE.
  13. Venugopalan, M., & Gupta, D. (2015, December). Sentiment classification for Hindi tweets in a constrained environment augmented using tweet specific features. In *International conference on mining intelligence and knowledge exploration* (pp. 664-670). Springer, Cham.
  14. Chiavetta, F., Bosco, G. L., & Pilato, G. (2016, April). A layered architecture for sentiment classification of products reviews in Italian language. In *International Conference on Web Information Systems and Technologies* (pp. 120-141). Springer, Cham.
  15. Alotaibi, S. S., & Anderson, C. W. (2016). Extending the knowledge of the Arabic sentiment classification using a foreign external lexical source. *Int. J. Nat. Lang. Comput*, 5(3), 1-11.
  16. O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks. arXiv preprint arXiv:1511.08458.