



# Categorizing and Labelling of Twitter Dataset using RCNN Model

Pradnya Saval<sup>1</sup>, Tahera Shaikh<sup>2</sup>

<sup>1,2</sup>Computer Engineering Department, Mumbai University

**Abstract:** *The fast advancement in the Web, web based business and interpersonal organizations realizes a lot of client created messages on the Web, for example, online surveys for items, administrations and web journals. Such messages as online audits are typically subjective and semantic arranged. The semantic introduction to separate and classify such messages appropriately is of incredible research and commonsense esteem. Order of writings is trying due of limited logical data and the scanty semantic data they regularly contain. The current research on content arrangement of writings fundamentally incorporates include based methods. The advancement of profound learning models has accomplished astounding outcomes in PC vision and discourse acknowledgment. In this content categorization method we will actualize the use of Convolution Neural Systems (CNNs) and Intermittent Neural Systems (RNNs) which acknowledges dynamic length input utilizing online networking information, for example, twitter. The proposed framework expects to build the execution measurements, for example, time and exactness to over 80% contrasted with other profound learning strategies.*

**Keywords:** *Neural network, Convolutional neural network, Recurrent neural network, Twitter dataset, Deep learning*

## I. INTRODUCTION

Content classification is the undertaking of appointing predefined classes to content records. It gives applied perspectives of reports and is utilized as a part of different applications in reality. For instance, cooking formulas are commonly sorted out by kind of fixings utilized, having a place with different land locales, formulas given by individual gourmet experts, scholastics can be classified to various spaces, therapeutic reports can be arranged to various tests et cetera. Spam separating is another application for content order where we channel or sort messages which are not to our significance. Content classification depends on words in which straightforward measurements of the mixes of words are performed till date. On the other hand, different analysts have discovered profound taking in valuable in isolating data from crude information running from PC vision to discourse acknowledgment to content arrangement. Characteristic dialect handling (NLP) has helped a great deal from the restoration of profound neural systems (DNN). There are different sorts of DNN designs in particular Recursive Neural System (RecursiveNN), Repetitive Neural System (RecurrentNN) and Convolutional Neural System (CNN). To beat the impediments of RNN calculation two new sorts of calculations have been created to be specific Long here and now memory (LSTM) and Gated Intermittent Unit (GRU)[1][4]. The following are the various models in deep learning:

### A. Recursive Neural Network (RecursiveNN):

RecursiveNN has turned out to be effective in building sentences. The RecursiveNN grab the clarification of a sentence utilizing a tree structure. The execution of this model wholly relies upon the portrayal of the tree structure giving the time multifaceted nature of  $O(n^2)$  where  $n$  is the length of the content. Utilizing this model could be tedious at the point when the model takes a shot at long sentences and even relationship between texts is difficult to speak to utilizing a tree structure. Along these lines it is expressed that demonstrating long sentences by recursive neural system is deficient [2]. The principle preferred standpoint of recursive neural organize is that it is proficient in building sentences anyway it is wasteful to demonstrate long sentences and reports. The execution of this model relies upon the literary tree structure and along these lines the execution corrupts as the tree structure increments relying on the content.

### B. Recurrent Neural Network (Recurrent NN)

Recurrent NN is a model that gives better reasonable data and is consecutive structures. It breaks down a word and stores the data of all the first words in a settled measured shrouded layer. This model is expressed as a one-sided display where all the following words in a sentence are successful than the past one because of which the data is accessible toward the finish of the report. Be that as it may, this can be an issue as the powerful data won't generally be available toward the finish of the report and the model uses its interior memory to process subjective arrangements of information sources. The time many-sided quality of RecurrentNN is  $O(n)$ .



RNN has the catch the logical data and can catch semantics of long messages as it can deal with self-assertive information and yield lengths. The utilization of RNN is for the most part utilized for content and discourse examination [3].

#### C. Convolutional Neural Network (CNN)

CNN reports higher execution contrasted with RecurrentNN. CNN processes the most educational n-grams of the considerable number of words in a sentence and concentrates the most profitable data. CNN can be utilized of estimation arrangement as the feelings are controlled by key expressions. It is an unprejudiced model and uses sustain forward ANN to reasonably segregate phrases utilizing a maximum pooling layer. In this manner CNN can catch semantic data superior to repetitive and recursive neural system. The time many-sided quality of CNN is likewise  $O(n)$ . Be that as it may, the past examinations tend that CNN utilize straightforward convolutional bits because of which it is hard to decide the measure of the window. . The utilization of CNN is for the most part utilized for picture and video handling [4].

#### D. Recurrennt Convolutional Neural Network (RCNN)

To defeat the negative marks of the above model RCNN was presented for content arrangement. To start with, the bi-directional intermittent structure is connected which acquaint less commotion contrasted and other customary window based neural system and subsequently it gives better relevant data. Second, we present a maximum pooling layer which decently figures out which include assumes a vital part contrasted with different words in the sentence. Consequently, the RCNN is a blend of repetitive neural system and convolutional neural system with a period multifaceted nature of  $O(n)$  [4].

## II. RELATED WORK

CNN is joined with LSTM without the technique for word division. The technique utilized was Burn CNN in view of Zhang's display where k-max pooling is actualized rather than max-pooling. Max-pooling is a technique for down testing where a sliding window is utilized on a column and chooses a cell with most extreme esteem and after that the window is passed to next layer.[5] K-max pooling on the other hand doesn't have a window rather it performs choosing activity for the whole column. The best k esteems with greatest esteem is chosen and go to the following layer. In this way, the creators have utilized Singe CNN technique by utilizing k-max pooling which has the ability to acknowledge any length of contribution before a completely associated layer and gives preferable exactness over other conventional word-level strategies. This technique can be connected to other common dialects utilizing word division for future research.

In [2], the creator has connected RCNN an intermittent structure to catch the relevant data utilizing a bi-directional repetitive structure alongside the maximum pooling layer which removes the most significant element and key parts in the content or report. Utilizing this one-sided display the later words are more overwhelming than the past one. In any case, the viability is decreased as the key segments won't generally be available toward the finish of the archive as the basic convolutional bits are hard to decide the window measure. Henceforth, it is imperative to decide the window measure for a basic convolutional portion to decide the highlights of a document.

In [3], the creators have given a methodical examination of the profound learning models, for example, CNN, LSTM and GRU utilized as a part of characteristic dialect preparing and expecting to choose the most suitable profound neural system models. The strategies actualized contains the preparation information utilizing fundamental setup without complex traps and have utilized for scanning parameters for each demonstrate independently. The examination specifies that the CNN show utilizes convolutional layer and GRU models the contribution from left to right and utilize the last layer as the after effect of the info. Because of this the variety in the bunch estimate and shrouded measure causes wavering which can be overcome in the event that it can be enhanced to expand the execution of the models.

The utilization of CNN display removes lexical and sentence level highlights. The model takes all the word tokens and changes to vectors utilizing word embeddings. The lexical highlights are then removed by the things introduce in the sentence. After the extraction both the techniques are joined to deliver the last component vector. To determine which combine of words should be marks to the separate word. The creator has utilized SemEval-2010 Errand 8 dataset which is uninhibitedly accessible and contains 10,717 examples alongside 8,000 and 2,717 preparing what's more, test occasions individually. [4].

In [5], the creator has utilized different extensive scale dataset, for example, AG's news with 4 classes, 120,000 what's more, 7,600 prepare and test tests individually alongside different datasets and actualized in the convolutional neural system. The examination of these dataset is among the different other conventional models, for example, pack of-words, n-grams and other profound learning models.

There are two kinds of CNN in particular, clear adjustment of CNN from picture to content and second is basic CNN utilizing bag-of- words. There different strategies used to arrange utilizing CNN will be CNN for picture, CNN for content, seq-CNN for content and bow-CNN for content. The correlation of mistake rates (%) appeared in the paper, demonstrates that seq-CNN beats different strategies for classification utilizing three datasets in particular Motion picture Audits (IMIDB) with 8.74, Hardware Item Surveys (Elec) with 7.78 and News Articles (RCV1) with 9.96 [6].

TABLE II  
RESEARCH TRENDS IN DEEP LEARNING

Problem	Methodology	Future Scope
Mix of CNN and LSTM for word division. [1]	Zhang based model uses Roast CNN where the maximum pooling idea used to utilize k elements.	word division in other dialect handling.
Utilizations one-sided demonstrate. Be that as it may, the viability is decreased because of confusion of segments toward the finish of the record. [2]	RCNN is utilized as the mix of utilizing the b-directional repetitive system alongside a maximum pooling layer which consequently gives the key parts in the document.	Determine the size bits in CNN.
Variety in group estimate and shrouded measure causes wavering. [3]	The strategies executed: Preparing information utilizing fundamental setup without complex traps. CNN utilizes convolutional layer then again GRU gets the contribution from left hand side and gets the outcome from last layer.	Optimization of concealed size and bunch size to expand the execution of the models.
Determination of which name to be allocated to match of words. [4]	Explore the connection and concentrate the lexical highlights. It likewise determines which highlight should be marked. In the paper SemEval-2010 Undertaking 8 dataset is used.	It physically picks the best list of capabilities and increment the execution.
Less accomplishment in the change of data recovery [5]	Using CNN on different dataset and make two Conv-Nets one expansive and other little and introduce the weight utilizing Gaussian distribution.	Apply CNN to a huge scope of dialect preparing undertaking.
1D structure of content information should misuse for precise anticipate work. [6]	Two strategies to be specific straight forward CNN and pack of-words change are applied.	Increasing the proficiency of CNN show with the assistance of sack of-n-grams approach.

### III. PROBLEM STATEMENT

In RNN the associations among units are as coordinated cycle that makes an inward state which enables the system to take a shot at dynamic conduct. Because of this RNN can characterize intermittent association in light of timestamps among different states. This issue in CNN models are utilized which utilizes convolutional layers called as convolutional layers. CNN can suit for content classification however are less noticeable in taking in the data from the information gave in the content. The CNN show is an unprejudiced model which utilizes the feed forward system and furthermore utilizes max-pooling layer rather than k-max pooling. This issue articulation is to beat every one of the defects said in the paper.

#### IV. PROPOSED SYSTEM

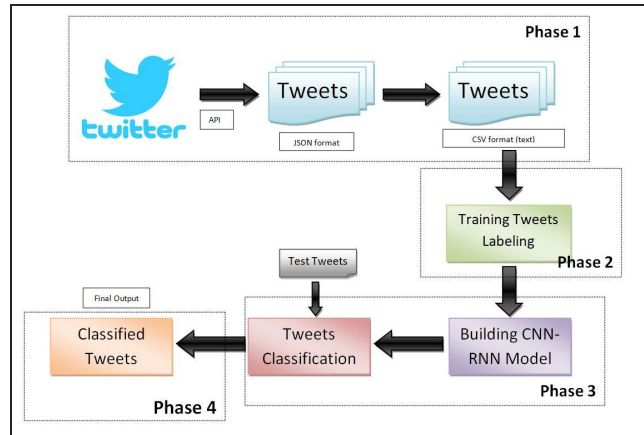


Fig. 1. Proposed System Architecture

In the proposed framework, is partitioned into different stages:

- 1) *Phase 1:* The tweets are downloaded from twitter Programming interface which is in .json arrange. The downloaded tweets should be changed over into .csv arrange by pre-handling it for additionally preparing of the dataset.
- 2) *Phase 2:* The changed tweets request to be prepared by marking it in particular infected, preventive, informative and others.
- 3) *Phase 3:* Once the information is prepared the dataset is given to the models and loss function is figured
- 4) *Phase 4:* The yield of the model is tried relying upon the loss function and the preparation gave with respect to the figure shown below.

The train model utilizes the dimension vector for each word and the dataset is split into training data and test data along with dev data. A directory is created which stores the data that is trained by the model along with the timestamp of that training instance. The model is then evaluated using the training data and the test samples i.e., the tweets. For every batch the data is trained and stored in the directory that is categorized as description and the category.

Once the data is trained the next step is to predict the models using the test data that are random flu tweets extracted from the Twitter API. The .csv file is loaded in the model along with the predicted and description as the heading. The predictions are then classified depending on the dictionary that is feed to the system and with the help of trained data the tweets are classified as preventive, infected and information.

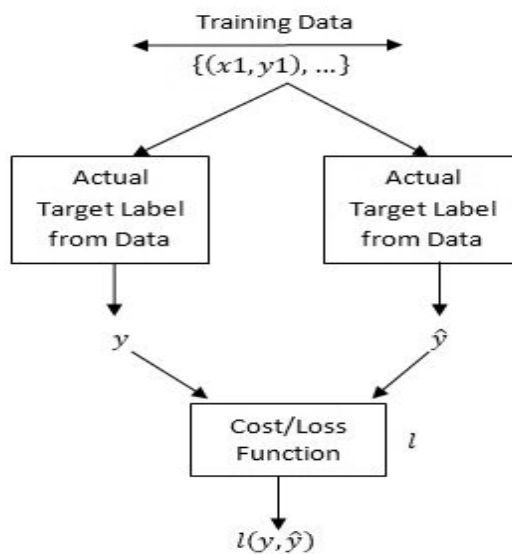


Fig 2. Loss function and of a neural network

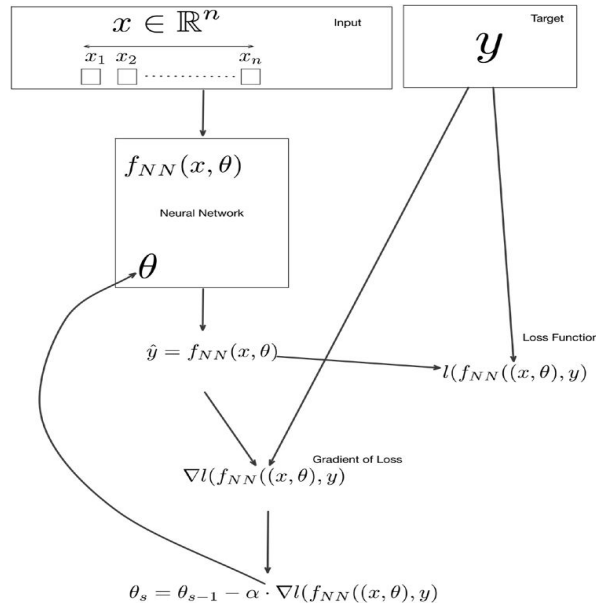


Fig 3. Flow of training data

The contribution to the model is given through preparing set. This information is given to the genuine target name and the neural system. The yield from the objective name and neural system is indicated as  $y$  and  $\hat{y}$  individually. Once the yield from both the squares have been gotten we process the forecast is right to what level i.e.,  $\hat{y}$  is contrasted and  $y$  and it is called as misfortune work. A misfortune work figures the disparity amongst  $y$  and  $\hat{y}$  is coordinated to as  $l$  as appeared in Fig 4. [9]

### V. DATASET

In the proposed strategy we utilize tweets from Twitter Programming interface. The dataset contains tweets of swine influenza or simply expressed as influenza which are refined from the removed tweets. An aggregate of 10,48,576 tweets are gathered similarly for five classes that is tainted, preventive, useful and others. The information is prepared utilizing the data above and the tweets are then named relying on the pre-characterized names gave amid preparing.

### VI. RESULTS AND DISCUSSION

The lead of the proposed CNN calculation is contrasted and other existing calculations. The outcome is appeared in Table V. Our anticipated model has a superior exactness than RNN and LSTM models. Be that as it may, the best model in this examination is dynamic CNN (DCNN), which creates more exactness than our model around 5% [1][3][16].

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	PREDICTED	Describe											
2	Infected	I got flu n coughed a lot. Now my voice is like monster voice. Rrr											
3	Infected	@lis bad head aching all rotten cough. Hasn't started making piggy noises yet tho!!!											
4	Infected	Barber just coughed on me in the chair. Pretty sure I now have swine flu											
5	Information	Research shows that a health care provider recommendation for a yearly #flu vaccine is very important to patients. #FightFlu.											
6	Information	t you can assess your patient vax needs & make a strong rec for them to get a flu vaccine.											
7	Prevention	A yearly #flu vaccine is the best way to prevent flu illness.											
8	Prevention	T everyday preventive actions & antivirals as recommended.											
9	Infected	i got flu and very cold cough											
10	Prevention	one need to have proper medicine and health care to prevent flu											
11	Infected	i am suffering from severe cough and cold											
12	Information	flu vaccine is very important to patients											
13	Prevention	get regular vaccination											
14	Prevention	God I am fed up of this flue since so many days											
15	Prevention	Take proper rest when you have flue											
16	Prevention	Take proper treatment while having flu											
17	Infected	i am suffering from cold											
18	Prevention	Drink lots of water and take rest											
19													

Fig 4. Final predicted result

TABLE III  
RESULTS WITH VARIOUS DATASETS

Method	20News	ACL	SST
LSTM	70.01	73.01	80.23
RNN	75.4	75.7	89.36
CNN	79.02	82.16	81.05
RCNN	82.83	90.24	80.36
Proposed - CNN	85.42	81.01	91.03

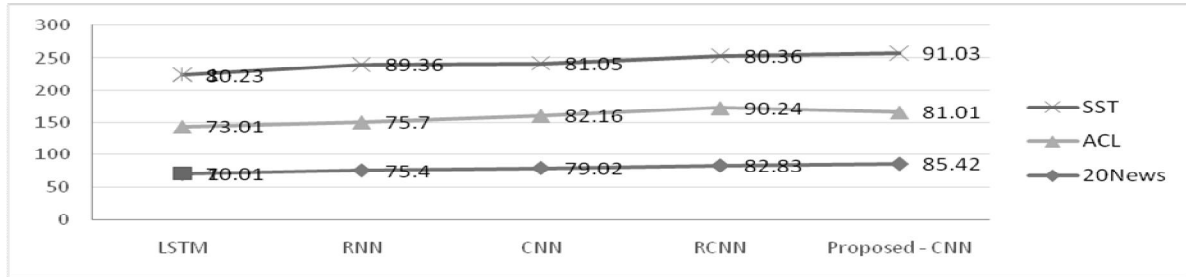


Fig 5. Comparison of various models with respect to different dataset [18] [19]

TABLE IV  
EXPERIMENTAL RESULTS

Method	Precision	Recall	F-measure	Accuracy
LSTM	73.01	73.1	73.01	88.23
RNN	75.4	76.5	75.7	89.36
Proposed - CNN	94.03	74.82	81.01	91.03

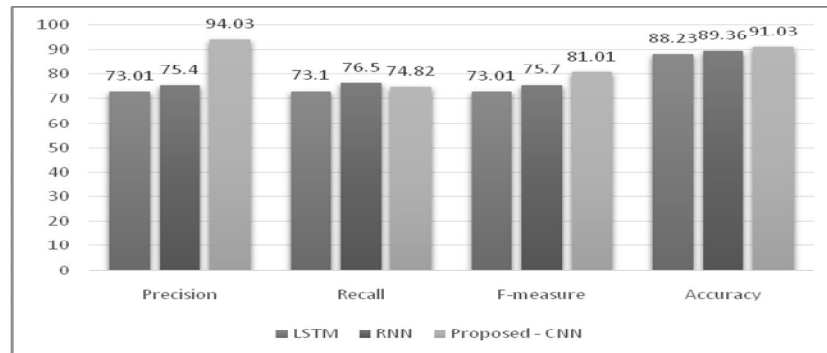


Fig 6. Final Comparison with respect to the parameters

We have inspected the model utilizing character input. The technique is directed by utilizing the mix of max and k-max pooling as said in the model depicted previously. The outcome in LSTM is awful contrasted with CNN as the criticism is likewise added to the model utilizing an overlook entryway which avoids the weights on the off chance that it is immaterial. CNN despite what might be expected is a fair model uses max-pooling and catches preferable semantic information over different models.

## VII. CONCLUSIONS

The outlined model is developed for proficient and compelling grouping to acknowledge variable length input. To perform classification errand in light of twitter information, the dataset is prepared utilizing Convolutions calculation. The CNN demonstrate contains input layer, gathering of units in a system named as widths, accumulation of widths (concealed layers) and a last layer. Our planned model which can acknowledge a more drawn out information gives precision in more superb way than a unique model with a settled length input. It likewise beats other neural system techniques e.g. LSTM, RNN aside from a word-level CNN. In future, CNN can be created to improve and adequately prepare the consecutive information for better execution.



## REFERENCES

- [1] Thanabhat Koomsubha, "A Character-level Convolutional Neural Network with Dynamic Input Length for Thai Text Categorization", 978-1-4673-9077-4/17/\$31.00 ©2017 IEEE.
- [2] Siwei Lai, Liheng Xu, Kang Liu, Jun Zhao, "Recurrent Convolutional Neural Networks for Text Classification", Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015.
- [3] Wengpeng Yin, Katharina Kann, Mo Yu, Hinrich Schütze, "Comparative Study of CNN and RNN for Natural Language Processing", Computation and Language, February 2017.
- [4] Daojian Zeng, Kang Liu, Siwei Lai, Guangyou Zhou, Jun Zhao, "Relation Classification via Convolutional Deep Neural Network", Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers, pages 2335–2344, Dublin, Ireland, August 23-29 2014.
- [5] Xiang Zhang, Junbo Zhao, Yann LeCun, "Character-level Convolutional Networks for Text Classification", Computation and Language, December 2015.
- [6] Rie Johnson, Tong Zhang, "Effective Use of Word Order for Text Categorization with Convolutional Neural Networks", Computation and Language, December 2014.
- [7] Tharani S, Dr. C. Yamini, "Classification using Convolutional Neural Network for Heart and Diabetics Datasets", International Journal of Advanced Research in Computer and Communication Engineering, ISO 3297:2007 Certified Vol. 5, Issue 12, December 2016.
- [8] Shiyao Wang, Zhidong Deng, "Tightly-coupled convolutional neural network with spatial-temporal memory for text classification", International Joint Conference on Neural Networks (IJCNN), 2017.
- [9] Rita Georgina Guimarães, Renata L. Rosa, Denise De Gaetano, "Age Groups Classification in Social Network Using Deep Learning", IEEE Access, Vol.5, 2017.
- [10] Nikhil Ketkar, Deep Learning with Python: A Hands-on Introduction, Apress, 2017.
- [11] Atsuya Oishia, \*, Genki Yagawa, "Computational mechanics enhanced by deep learning", 0045-7818 © 2017 Elsevier Comput. Methods Appl. Mech. Engrg. 327 (2017) 327–351
- [12] Carlos Affonso a , \*, André Luis De Biaso Rossi, "Deep learning for biological image classification", Expert Systems With Applications, 0957-4174/© 2017 Elsevier, 114-122
- [13] Bilal Jan, "Deep learning in big data Analytics: A comparative study", Computers and Electrical Engineering, pp.1–13 , Elsevier, 2017
- [14] Geert Litjens, "A survey on deep learning in medical image analysis", Medical Image Analysis, pp.60–88, Elsevier, 2017
- [15] Y. Kim, "Convolutional neural networks for sentence classification," in Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2014, pp. 1746–1751.
- [16] X. Zhang, J. Zhao, and Y. LeCun, "Character-level convolutional networks for text classification," in Advances in Neural Information Processing Systems, 2015, pp. 649–657.
- [17] P. Vateekul and T. Koomsubha, "A study of sentiment analysis using deeplearning techniques on Thai Twitter data," in Computer Science and Software Engineering (JCSSE), 2016 13th International Joint Conference on, 2016.
- [18] N. Kalchbrenner, E. Grefenstette, and P. Blunsom, "A convolutional neural network for modelling sentences," in Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, {ACL} 2014, 2014, pp. 655–665.
- [19] C. N. dos Santos and M. Gatti, "Deep Convolutional Neural Networks for Sentiment Analysis of Short Texts," in COLING, 2014, pp. 69–78.