

Chapter 35

Comparison of Traditional Machine Learning and Deep Learning Approaches for Sentiment Analysis



Dhvani Kansara and Vinaya Sawant

1 Introduction

Natural language processing (NLP) is the field lying at the intersection of computer science, artificial intelligence and linguistics. Sentiment Analysis, which is a sub-domain of NLP, is a process of determining the emotional tone of a comment based on words contained in it, in order to understand the attitude and opinion behind it. This process is also called opinion mining. It has multiple applications mainly in the form of social media monitoring. In recent years, there have been great advancements in this field. Various machine learning algorithms have shown to prove effective in categorizing the sentiment from a text. Recent advances have been focusing on using deep learning algorithms for this purpose. This proliferation is due to the fact that opinions are central to almost all human activities and are key influencers of our behaviors. Our beliefs and perceptions of reality, and choices we make, are, to a considerable degree, conditioned upon how others see and evaluate the world [1].

In this paper, we discuss two paradigms: traditional approaches for classification which have been in use since the past few decades and the recent breakthroughs leveraging deep learning algorithms. The models work on raw data further cleaned and pre-processed, considering removal of stop words, punctuations and mark-ups with stemming (process of reducing inflected words to their word stem) and then checking the overall sentiment by assigning polarity based on resultant cleaned sentence. The models also take into account inversion terms or negations (such as—“not bad”) which reverse the polarity of the sentence as a whole. Later, this data is fed into a learning-based model that uses a supervised learning algorithm. In case of

D. Kansara (✉) · V. Sawant (✉)
Dwarkadas J. Sanghvi College of Engineering, Mumbai University, Mumbai, India
e-mail: dhvani.djk@gmail.com

V. Sawant
e-mail: vinaya.sawant@djsce.ac.in

© Springer Nature Singapore Pte Ltd. 2020
H. Vasudevan et al. (eds.), *Advanced Computing Technologies and Applications*, Algorithms for Intelligent Systems,
https://doi.org/10.1007/978-981-15-3242-9_35

traditional models, this text is pre-processed using the bag of words methodology. It gives a combination of types of semantic features that attempt to model the syntactic structure of sentences, intensification, negation, subjectivity or irony [2]. In this paper, we use Naive Bayes, Logistic Regression and Random Forest as traditional classification approaches. On the other hand, deep learning models learn from multiple layers of representations or features of data and produces results accordingly. For this, we use vector representations of sentiments as inputs to the deep layers. We use modified Recurrent Neural Networks (RNN) or Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNN) and a combination of CNN and LSTM to obtain experimental results. We then analyze the accuracies obtained in all these algorithms.

We present each paradigm along with the description of pre-processing phases and an overview of the algorithms used. We also try to analyze reason behind difference between the accuracies. We use three different datasets for this purpose and classify movie reviews, hotel reviews and tweets, thereby concluding that deep learning models prove to have an upper edge in all the cases.

2 Methodology

2.1 *The Traditional Approach*

Extensive research has been carried out in the past few years to apply basic machine learning algorithms for sentiment analysis. We leverage this research to ensure that our approach will provide the best results. In order to obtain a good performance accuracy, the pre-processing phase is carried out prior to the classification process.

Pre-processing. The reviews may be taken directly from a Web site, so it may include html mark-up. Thus, we first use the BeautifulSoup package to strip all html mark-up. We then strip the stop words using the Porter Stemming package in Natural Language Toolkit (NLTK) to stem and remove stop words like “a,” “and,” “is,” “the” because these are frequently occurring words which do not carry any significant meaning. Then, we convert this categorical data of cleaned reviews and tweets into numeric data so that we can apply machine learning algorithms on them. We use an approach called bag of words [3] for this. This model creates a dictionary by taking all words in the dataset and then assigns integers according to the count of each word of the dictionary appearing in the given text. For example, consider the following two sentences:

Sentence 1: “The movie had amazing actors”

Sentence 2: “The view from the hotel room was amazing”

So, our vocabulary will be as follows:

{the, movie, had, amazing, actors, view, from, hotel, room, was}.

To get our bag of words, we count the number of times the word from the vocabulary list occurs in each of our sentences.

In Sentence 1, “the,” “movie,” “had,” “amazing,” “actors” each appears once, so the feature vector for Sentence 1 is: {1, 1, 1, 1, 1, 0, 0, 0, 0, 0}.

Similarly, the feature vector for Sentence 2 is: {2, 0, 0, 1, 0, 1, 1, 1, 1, 1} In this way, we obtain a sequence of integers for each review and tweet.

In the datasets, we have a very large number of reviews, which will give us a large vocabulary. To limit the size of the feature vectors, we chose a maximum vocabulary size of 5000 words, i.e., we consider the top 5000 most frequently occurring words in our vocabulary. (Considering we already removed stop words, this is enough.) Thus, we obtain a list of size 5000 with integer values mapped according to the vocabulary for each row.

The Algorithms Part I. The dataset is divided into 80% tuples as training data and remaining 20% tuples as test data and the algorithms are applied on it. Accuracy is calculated based on how correctly the model predicts the test data. Three traditional classification algorithms were experimented,

- i. Naive Bayes
- ii. Logistic Regression
- iii. Random Forests.

Naïve Bayes Algorithm. It is based on the formula that, for a label y , we have the independent data feature x_i as in Eq. (1).

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y) \quad (1)$$

where $P(y|x_1, \dots, x_n)$ is the posterior probability stating that x belong to class y . To create the classifier model, the probability of given set of inputs for all possible values of the class variable y is found, the output with maximum probability is the final classification result. This can be expressed mathematically as in Eq. (2).

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y) \quad (2)$$

This algorithm assumes the property of conditional independence among all attributes, which may not be true in all cases thus reducing the accuracy to a certain extent [4].

Logistic Regression. It is one of the most commonly used binary classification algorithm in machine learning which gives discrete binary output between 0 and 1 [16]. It uses a logistic function, also called as sigmoid function whose equation is as shown in Eq. (3)

$$f(x) = \frac{1}{1 + e^{-x}} \quad (3)$$

Taking into consideration the weight of each input variable the output is predicted, this weighted input is passed to the sigmoid function to obtain a probability value between 0 and 1. For any real number input, the sigmoid function maps it to a value between 0 and 1. These values will then be transformed into either 0 or 1 using a threshold classifier.

Since it is a discriminative classifier which does not rely on the assumption of attribute independency as in case of Naive Bayes, thus this model may generally have a better accuracy comparatively.

Random Forest Classifier. It is formed by a collection of multiple decision tree classifiers collectively called as a “forest.” Decision tree classifier concept is based on the rule-based system [17]. Individual decision trees are generated using selection of attribute at each node to determine the split. This node is selected in a strategic manner such that the most efficient tree is formed. The flow of the tree to classify a certain tuple will be in the form of an if-then rule. The same set rules can be used to perform the prediction on the test dataset. Multiple such trees are formed with different combinations of training and testing data. Then, it takes the test features and uses the rules of each randomly created decision tree to predict the outcome and stores the predicted outcome. Votes for each predicted result are calculated. Result which has received the maximum votes is considered as the final prediction. In this, multiple trees (different rules) are generated from given input instances or attributes and majority of the votes of the classifier output is considered as the final class [5]. Overfitting is avoided in this algorithm because every time a random set of inputs are used to build the decision tree. We use 100 estimators, i.e., 100 different decision trees to train the model and “entropy” as splitting criteria for determining the best split point.

2.2 Deep Learning Approach

Deep learning architecture consists of input layers, hidden layers and output layers. Each layer is made up of neurons, which are basically weight matrices. The higher the weight, the more influence layer has over the next. The input is passed through the hidden layers to give a certain output. Backpropagation is used to minimize the loss incurred between predicted and actual outcome. It does so by changing the weight matrices in the dense layers such that error is minimized. In this way, the learning occurs.

Pre-processing. The cleaning of the datasets is done as in the previous case, i.e., by removing html mark-up, punctuations and stop words. For input to the dense layers, we form a vector in latent space for each review, and this approach is called Word-to-Vectors. We use tokenizer for dividing sentences into tokens (or words). Unique words are identified and unified in a list, from all inputs taken into consideration collectively. We then assign an integer value for each of the unique word in the list. So now we have a dictionary of words each assigned a specific integer number. Then,

idx ▲	Type	Size	Value
0	unicode	1	stuff go moment mj start listen music watch odd documentari watch wiz ...

Fig. 1 Cleaned review

idx ▲	Type	Size	Value
0	list	219	[458, 24, 166, 6988, 85, 878, 83, 11, 793, 499, ...]

Fig. 2 Token representation of review

each review is represented by a list of these integers which form a part of that review. The transformation is as shown in Figs. 1 and 2.

We pad the short reviews with 0s in the beginning and truncate the longer reviews. During training, the model will learn that the 0s carry no information thus preserving the content in each case and equalizing the length of each input, which is necessary for computation. We constrict a total size of 300 for each review.

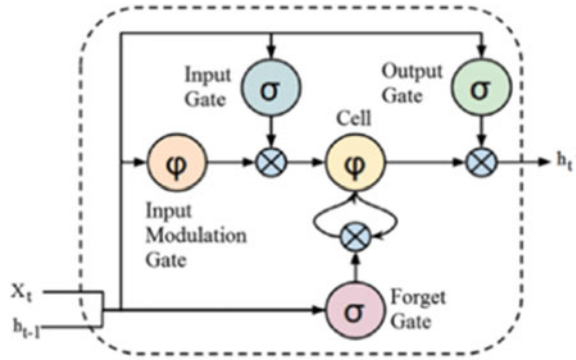
The Algorithms Part II. The datasets are divided into 80% tuples as training data and remaining 20% tuples as test data and the algorithms are applied on it. Accuracy is calculated based on how correctly the model predicts the test data. The following algorithms were experimented:

- i. Recurrent Neural Network (Long Short-Term Memory)
- ii. Convolutional Neural Network
- iii. Convolutional Neural Network + Long Short Term Memory.

Recurrent Neural Network (Long Short-Term Memory). RNNs remember all the relations while training itself and are thus used particularly for sequential data [6]. The output of the previous state also goes into the input of the next state along with current input vector, this helps the model to remember context while training. That is why it is used for sentiment analysis and due to its nature it can be used to correctly predict sentiments of statements such as “Though initially I liked the movie, but toward the end it became dull and boring” this sentence has a negative sentiment but also positive words like “liked” which can confuse the system, but not in case of recurrent neural network, which can remember long-term dependencies.

In traditional recurrent neural network, during the backpropagation phase, if the values of weight matrix of the recurrent neural layers are very small (less than 1.0) that ultimately they do not have any effect on the input signal and thereby preventing the neural network from training further. The problem is called vanishing gradient problem because the gradient signal vanishes slowly when passing through different layers [7]. Thus, RNNs have problems learning long-range dependencies, i.e., relations between words that are several steps apart. Conversely, we could also face exploding gradient problem if the weight matrix values become too huge, causing training to diverge. For this reason, we use a modified version of recurrent neural network called Long Short-Term Memory (LSTM) Neural Network which is capable

Fig. 3 LSTM functioning



of learning long-term dependencies. We need to update information less chaotically and account for all the learning and the memory that it has learnt over a vast period of time, in order to make more accurate [15].

Labels in Fig. 3 and equations below represent:

X_t –current input vector | b –coefficients | h_{t-1} –Output of previous block | W —weights | h_t —output of current block | σ —sigmoid function | φ —tanh function.

Figure 3 [8] describes the functioning of a LSTM model. The cell state is the most important component which “remembers” values over arbitrary time periods thus giving LSTM a long-term memory; it represents all the learnings across time.

An LSTM cell consists of three gates as shown in Fig. 3:

1. A forget gate which is used to formulate an attention mechanism which helps filter out relevant data by knowing what to forget and what to remember. It takes current input(X_t) and previous output(h_{t-1}) and passes it through a sigmoid function and judges if the data is worth remembering, which outputs 1 if data is to be remembered or 0 if it is to be forgotten. Thus, its equation is given by Eq. (4).

$$f_t = (W_f[h_{t-1}, X_t] + b_f) \tag{4}$$

2. An input gate decides what new information we are going to store in the cell state by involving a sigmoid function to X_t and h_{t-1} . The equation of input gate is given by Eq. (5).

$$i_t = (W_i[h_{t-1}, X_t] + b_i) \tag{5}$$

The tanh layer (gives output between -1 and $+1$) is responsible for creation of vector of new candidate value, \tilde{C}_t as given by Eq. (6).

$$\tilde{C}_t = \tanh(W_C[h_{t-1}, X_t] + b_C) \quad (6)$$

3. This new candidate value is multiplied with the input gate and forget gate output along with previous cell state, C_{t-1} is used to formulate the new cell state C_t as given in Eq. (7).

$$C_t = (f_t * C_{t-1}) + (i_t * \tilde{C}_t) \quad (7)$$

An output gate decides what information we are going to output as in Eq. (8).

$$o_t = (W_o[h_{t-1}, X_t] + b_o) \quad (8)$$

Thus, output of current block is given by multiplying output gate value with new cell state passed into a tanh layer as shown in Eq. (9) [7–9].

$$h_t = o_t * \tanh(C_t) \quad (9)$$

Thus, this model is able to remember context by storing relevant information in its long-term memory.

Embedding Layer: Above Fig. 4 shows layout of our model. Before the LSTM layer, we add an embedding layer whose role is to map semantic meaning to geometric space. This is done by associating a numeric vector to every row in the dataset, such that the distance between any two vectors would capture part of the semantic relationship between the two associated rows. For example, consider two words “potato” and “panda” which are not semantically related so their vectors would be far apart but “green” and “grass” will have closer values. The output values will be in the form of vectors with related word vectors having close by distances, thus closer values. Also, vector arithmetic is applicable to these word embeddings, for example,

$$\text{Vectorman} - \text{Vectorwoman} \approx \text{Vectorking} - \text{Vectorqueen}$$

Thus, vectors are formed based on the contexts of the words. So for example, if the machine assigns -1 for gender male and $+1$ for female, eventually learning these will enable queen getting a gender dimension of $+0.97$ and king of -0.95 and for say for mango and peach sort of genderless. Other dimensions may include whether or not is it a fruit, ages, cost, is it a noun or verb, size, loyalty, etc., thus one can come up with multiple dimensions to represent an item. In our case, vectors will be formed based on how closely related two reviews (rows in the form of tokens) are. We use



Fig. 4 Layers in training LSTM model

an embedding dimension of 100 which means that each review that is in our dataset is mapped into a 100-dimension dense vector of floating numbers. The vector will look like- (0.15,... 0.23,... -0.55).

This output is fed into the input of LSTM layer, in this layer, we apply a dropout of 20% to reduce overfitting, which drops 20% neurons in this layer. Another hidden layer is added following this which uses the sigmoid function to give a binary output.

Convolutional Neural Networks. Convolutional neural networks were designed to honor the spatial structure in image data while being robust to the position and orientation of learned objects in the scene [10]. They use pixel values and feature matrix for this. The same principle can be used to help learn structure of paragraphs of words, namely the techniques invariance to the specific position of features. This is applied over one-dimensional sequence of reviews to learn its sequence in a similar fashion. So, instead of image pixels, the inputs to CNN model are sentences represented as matrix. Each row corresponds to one word, these are word embeddings, explained previously. The width of the feature filter is same as the width of the input matrix and height is varied.

The input sentence is passed into a feature detector or a filter, which are basically vectors used to learn the sequence of the statements. As the filter is sliding (or convolving), it is multiplying its weight values with the original word embedding value matrix of the review. The multiplications are summed up to a single number, which is a representative of the receptive field. After scanning, we get an array which is called as the feature map. To scan the input, we make use of 100 filters each of size 3, which forms the first convolutional layer, which are responsible for feature selection. We also used rectifier activation function to remove the linearity (if any) present in the data. Following this layer is the max-pooling layer which is used to reduce the spatial size of the representation, which in turn reduces the computation complexity of the network, this reduces the number of dimensions. Here, care is taken that the most important information is retained. The output from here is then flattened to form a 1D matrix which is fed into input of the next dense hidden layer, which uses rectifier activation function that transforms the incoming signals of the neurons and pass these signals to the input of the next hidden layer, which uses sigmoid as its activation function to classify a binary output. CNN has a special spatially local correlation by enforcing a local connectivity pattern between neurons of adjacent layers. Such a characteristic is useful for classification in sentiment analysis where we expect to find a certain pattern of sequences that may be present in different parts of a sentence [11, 12]. Figure 5 shows different layers used in training the CNN model.

Convolutional Neural Network + Long Short Term Memory. This model is a hybrid of the previous two models. The CNN will gather distinct features for positive and



Fig. 5 Layers in training CNN model



Fig. 6 Layers in training CNN + LSTM model

negative sentiment and these spatial features will then be learned as sequences by the LSTM layer [13].

Accordingly, we use a one-dimensional CNN and max-pooling layer after the Embedding layer which is then fed as the consolidated features to the LSTM. We use a set of 100 filters with a filter length of 3. The pooling layer used is of the standard length 2 to halve the feature map size. The max-pooled output is then fed as the input to the LSTM layer, followed by a hidden layer, which uses sigmoid as activation function to classify a binary output. Figure 6 shows layers in the model.

3 Experiments

The following three datasets were used:

1. IMDB Movie review dataset

This dataset was taken from the Kaggle Web site. It consists of 25,000 rows each containing three attributes-A unique identifier, review for a movie and its corresponding binary sentiment-1(positive) and 0(negative).

2. Twitter Dataset

This dataset was also taken from the Kaggle Web site. It consists of three columns, tweet id, tweet text and its corresponding binary sentiment. Since we are considering only binary outcomes, i.e., 0 or 1, the accuracies obtained in testing this dataset are low comparatively because we are not considering a neutral output, which may be frequent among tweets.

3. Datafiniti's Business Database

This dataset was taken from data.world Datasets. It contains data for 1000 hotels with 35,000 rows and following attributes: hotel location, name, rating, review data and review username. From this we use hotel review data and corresponding rating. Since the rating is on a scale of 5, we consider rating >2 as positive, i.e., assign binary 1 sentiment and rating <3 as negative, i.e., assign binary 0 sentiment during the pre-processing phase.

In all three cases, the models learn from the training dataset about the words which help classify a review as positive or negative. Then, it predicts the outcome (sentiment) for reviews in the test dataset. Thus, the input is the given review/tweet and the output is predicted label- 0 or 1 of that input text. To find the accuracy, we

compare the predicted labels with the actual labels of those corresponding reviews in the test dataset. Accuracy is therefore given by the formula in Eq. (10).

$$\text{Accuracy} = \frac{\text{Correctly classified rows}}{\text{Total number of rows}} \quad (10)$$

In the deep learning approach, for training we use 1 epoch with a batch size of 32 in each case, meaning we update the weights once in each hidden layer, with 32 training tuples at a time, iteratively, till all training tuples exhaust. Before training, we configure the learning process by compiling with “adam” as optimizer function for efficiently calculating the gradient descent during the backpropagation for minimizing the loss incurred during while updating weights. To calculate this loss, we use binary_crossentropy as the loss function for all our models.

We also ensured that we chose the best methods of pre-processing in each case such that accuracy is not compromised.

4 Results

All datasets were experimented and their accuracies are as obtained in Table 1. From this, we can infer that in all three cases, deep learning algorithms prove superior over traditional machine learning algorithms. The drastic characterization of change in accuracies can be seen in Fig. 7.

From the above results, it is evident that deep learning algorithms have an edge over traditional machine learning algorithms. Elaborately because, a deep learning technique learns categories incrementally through its hidden layer architecture, defining low-level categories like letters first then little higher level categories like words and then higher level categories like sentences. Each node in the hidden layer is given a weight that represents the strength of its relationship with the output and as the model develops, the weights are adjusted.

Table 1 Accuracies (in%) of each classification algorithm on three datasets

Dataset Algorithm	IMDB reviews	Hotel reviews	Twitter tweets
Naïve Bayes	71.98	70.9	61.28
Logistic Regression	85.48	80.12	72.90
Random Forests	84.98	80.37	72.44
LSTM	87.58	81.29	74.54
CNN	88.98	81.28	74.44
CNN + LSTM	88.28	81.19	74.92

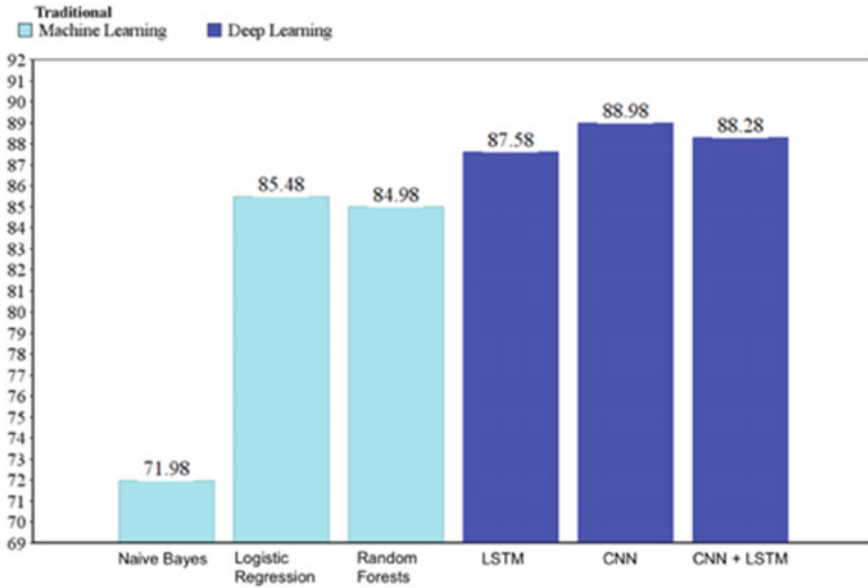


Fig. 7 IMDB movie reviews accuracies comparison for machine learning and deep learning approaches

The mascot which gives this advantage to deep learning models in NLP is the *embedding layer*, the fact that featured representation is possible using embeddings is what outsmarts them compared to traditional machine learning algorithms. Deep learning ensures that they understand the language that they are interpreting by taking word embeddings as input which retain context information and pass them to the hidden layers to learn these features during the training phase itself and thereby making more accurate predictions. Also, the ML algorithms cannot learn the sequence of words which form the sentences, thus making them less suitable for sentiment analysis tasks.

Additionally, in traditional machine learning techniques, most of the applied features need to be identified by a domain expert in order to reduce the complexity of the data and make patterns more visible for learning algorithms to work. We need to perform dimensionality reduction to find the best features to pass over the ML algorithms, which is not the case in deep learning. The biggest advantage of deep learning algorithms is that they try to learn high-level features from data in an incremental manner thus giving better performance right off the bat. This eliminates the need of domain expertise and hardcore feature extraction [14]. Manually designed features are often over-specified, incomplete and take a long time to validate whereas learned features are easy to adapt.

Though for smaller datasets traditional ML algorithms may outperform deep learning algorithms, deep learning requires sufficiently large amounts of data to work well. Also, CPU utilization is about 10 times more in case deep learning models, thus increasing the time employment.

5 Conclusion

The classification performance of the different models on movie reviews, hotel reviews and tweets gave a rough idea of the utility of these models for sentiment analysis. The performance of the deep learning models was overwhelming. Major reason for this being that the hidden layers understand the copious amount of data exceedingly well and have a deeper semantic understanding of the sequences, plus contextual relationship is retained using word embeddings which boosts the understanding of the model. In fact, they promise to perform much better as seen by the improved accuracies in the above experiments. The complementary advantage is that we can use deep learning models when there is lack of domain understanding for feature introspection, because one has to worry less about feature engineering. Thus, owing to profound contextual understanding deep learning models outperform traditional machine learning algorithms for most sentiment analysis and other NLP-based tasks.

References

1. Zhang L, Wang S, Liu BF Deep learning for sentiment analysis: a survey
2. Maite Taboada F Sentiment analysis: an overview from Linguistics. *Art Annual Rev Linguistics*. <https://doi.org/10.1146/annurev-linguistics-011415-04051>
3. <https://www.kaggle.com/c/word2vec-nlp-tutorial#what-is-deep-learning>. Last Accessed 30 July 2018
4. Rish I An empirical study of the naive Bayes classifier. T.J. Watson Research Center
5. Jiawei H, Micheline K, Jian P (2012) Data mining concepts and techniques. The Morgan Kaufmann series in data management systems. Morgan Kaufmann Publishers, Elsevier
6. Arras L, Montavon G, Müller KR, Samek WF Explaining recurrent neural network predictions in sentiment analysis
7. <http://deeplearning.net/tutorial/lstm.html>. Last Accessed 30 July 2018
8. <https://medium.com/@kangeugine/long-short-term-memory-lstm-concept-cb3283934359>. Last Accessed 30 July 2018
9. <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>. Last Accessed 30 July 2018
10. Hochreiter S, München, Germany, Schmidhuber J, Switzerland F Long short term memory. *Neural Comput Arch* 9(8)
11. Zhang Y, Wallace BCF A Sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. University of Texas, Austin
12. <http://www.joshuakim.io/understanding-how-convolutional-neural-network-cnn-perform-text-classification-with-word-embeddings/>. Last Accessed on 02 Aug 2018
13. <https://machinelearningmastery.com/sequence-classification-lstm-recurrent-neural-networks-python-keras/>. Last Accessed on 30 July 2018

14. Singhal P, Bhattacharyya PF Sentiment analysis and deep learning: a survey. Singhal P, Bhattacharyya P (eds). Indian Institute of Technology, Powai
15. Timmaraju A, Khanna VF Sentiment analysis on movie reviews using recursive and recurrent neural network architectures. Stanford University
16. Cramer JJ The origins of logistic regression. Tinbergen Institute discussion paper
17. <https://medium.com/@williamkoehrsen/random-forest-simple-explanation-377895a60d2d>