Text based Sentiment Analysis using LSTM

Dr. G. S. N. Murthy, Shanmukha Rao Allu, Bhargavi Andhavarapu, Mounika Bagadi, Mounika Belusonti Department of Computer Science and Engineering Aditya Institute of Technology and Management Srikakulam, Andhra Pradesh

Abstract— Analyzing the big textual information manually is tougher and time-consuming. Sentiment analysis is a automated process that uses computing (AI) to spot positive and negative opinions from the text. Sentiment analysis is widely used for getting insights from social media comments, survey responses, and merchandise reviews to create data-driven decisions. Sentiment analysis systems are accustomed to add up to the unstructured text by automating business processes and saving hours of manual processing. In recent years, Deep Learning (DL) has garnered increasing attention within the industry and academic world for its high performance in various domains. Today, Recurrent Neural Network (RNN) and Convolutional Neural Network (CNN) are the foremost popular types of DL architectures used. We do sentiment analysis on text reviews by using Long Short-Term Memory (LSTM). Recently, thanks to their ability to handle large amounts of knowledge, neural networks have achieved a good success on sentiment classification. Especially long STM networks.

Keywords—Sentiment Analysis, Text Classification, LSTM, Deep Learning

I. INTRODUCTION

Sentiment analysis is that the computerized process of the higher cognitive process to an opinion a couple of given subjects from a transcription. in an exceedingly present generation, we create quite 1.5 quintillion bytes of information daily, sentiment analysis has become a key tool for creating a sense of that data. it absolutely was utilized by the businesses to induce key insights and automate every kind of process for their business development. Sentiment Analysis [1] is also called opinion mining. Sentiment analysis isn't only a sentiment mining but also contextual mining of text which identifies and extracts subjective information in source material and helping a business to know the social sentiment of their service, brand or product while monitoring online conversations. Sentiment Analysis is that the most used text classification tool that analyses an incoming message and tells whether the essential opinion is positive or negative. Sentiment analysis will be applied at different levels of scope like Document-level sentiment analysis obtains the sentiment of an entire document or paragraph. Sentence level sentiment analysis obtains the results of one sentence. Sub-sentence level sentiment analysis obtains the results of sub-expressions within a sentence.

A. Why sentiment analysis is important?

It's estimated that 80% of the world's data is unstructured and not organized during a pre-defined manner. Most of this comes from text data, like reviews, emails, chats, social media, surveys and articles. These texts are usually difficult and time-consuming to investigate and understand. The sentiment analysis system authorizes company to create sense of this huge amount of unstructured text by automating business processes, saving hours of manual processing [2] and getting actionable insights.

Recurrent Neural Networks (RNNs) are one of the most prevalent architectures because of the ability to handle variable-length texts. Humans can't analyze from scratch every second. Any human can understand each word based on his understanding of previous words. He doesn't throw everything away and start thinking from scratch again. His thoughts have persistence. Traditional neural networks can't do this, and it seems like a speed process is coming. For example, imagine a human want to classify what kind of event is happening at every point in a movie. It's not clear how a traditional neural network could use its reasoning about previous events in the film to inform later ones. Recurrent neural network addresses can face this type of issues. They are networks with multiple loops in them, allowing information to continue. Though RNNs are capable of modeling long sequential data theoretically they fail to represent long sequences in real time applications [3].

Recently, LSTM is most popular to deal with sentiment classification. LSTM is proposed by Hoch Reiter and Schmid Huber in 1997 and was refined and popularized by many people in the following work. They work tremendously well on large different types of problems and are now widely used. LSTMs are explicitly designed to ignore the long-term dependency problem [4]. Remembering information for a long time is practically their default behavior, not something they struggle to learn. All recurrent neural networks have the form of a chain of repeating modules of the neural networks. In the level of RNNs, this repeating module having a very simple structure, such as a single tanh layer. The IMDB benchmark dataset is used for our experimental studies that contain movie reviews that are classified as being positive or negative.

An Example for positive and negative words

II. PROPOSED WORK

Long short-term memory (LSTM) is a synthetic recurrent neural network (RNN) architecture employed in the sphere of deep learning. Unlike standard feedforward neural networks, LSTM has feedback connections.

LSTM [5] networks are well-suited to classifying, processing, and making predictions supported statistic data since there may be lags of unknown duration between important events in a very statistic. LSTMs were developed to accommodate the exploding and vanishing gradient problems that may be encountered when training traditional RNNs. Relative insensitivity to gap length is a bonus of LSTM over RNNs, hidden Markov models, and other sequence learning methods in numerous applications. There are several architectures of LSTM units. a typical architecture consists of a cell (the memory a part of the LSTM unit) and three "regulators", usually called gates, of the flow of knowledge inside the LSTM unit: an input gate, an output gate and a forget gate. Some variations of the LSTM unit don't have one or more of those gates or even produce other gates. as an example, gated recurrent units (GRUs) don't have an output gate.



Max Sequence Length

LSTM with a forget gate

The compact forms of the equations for the forward pass of an LSTM unit with a forget gate are:

$$egin{aligned} f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \ i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \ o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \ ilde{c}_t &= \sigma_h(W_c x_t + U_c h_{t-1} + b_c) \ c_t &= f_t \circ c_{t-1} + i_t \circ ilde{c}_t \ h_t &= o_t \circ \sigma_h(c_t) \end{aligned}$$

where the initial values are $c_0=0$ and $h_0=0$ and the operator ^o denotes the Hadamard product (element-wise product). The subscript *t* indexes the time step. In this model, σ is the sigmoid activation function, tanh the hyperbolic tangent activation function, Xt the input at time t, W_i, W_c, W_f, W_o, U_i, U_c, U_f, U_o are weight matrices to regulate the input and b_i, b_c, b_f, b_o are bias vectors.

III. ARCHITECTURE OF PROPOSED NETWORK USED



A. Raw Text

We are using IMDB movies review [6] and Amazon Product datasets used to train and validate our models. In total, these datasets contain tweets labeled as either positive or negative. If it is stored in your machine in a text file then we just load it. Then we convert the text to lower case and remove punctuation. We have got all the strings in one huge string. Now we have to separate out individual reviews and store them in individual list elements. Like, [review_1, review_2, review 3...... review n].

B. Tokenizer

Tokenization is that the process of tokenizing or splitting a string, text into an inventory of tokens. One can consider token as parts sort of a word could be a token in a very sentence, and a sentence could be a token in a very paragraph. Tokenizing (splitting a string into its desired constituent parts) is key to all or any NLP tasks. there are no single right thanks to doing tokenization. the correct algorithm depends on the appliance. I suspect that tokenization is even more important in sentiment analysis than it is in other areas of NLP, because sentiment information is often sparsely and unusually represented a single cluster of punctuation like >:-(might tell the whole story.

- 1. Create Vocab to Int mapping dictionary
 - In most of the NLP tasks, you will create an index mapping dictionary in such a way that your frequently

occurring words are assigned lower indexes. One of the most common ways of doing this is to use the Counter method from the Collections library.

- 2. Encode the words
- So far, we have created a list of reviews and index mapping dictionaries using vocab from all our reviews. All this was to create an encoding of reviews (replace words in our reviews by integers) what we have created now is a list of lists. Each individual review is a list of integer or floated values and all of them are stored in one huge list.
- 3. Encode the labels
- This is simple because we only have 2 output labels. So, we will just label 'positive' as 1 and 'negative' as 0.
- This class allows vectorizing a text corpus, by turning each text into either a sequence of integers (each integer being the index of a token in a dictionary) or into a vector where the coefficient for each token could be binary, based on word count, based on term frequency-inverse document frequency.

C. Embedding

Word Embedding [7] emerged from the field of Natural Language Processing (NLP) which is an intersection of Computer Science, Artificial Intelligence, Machine Learning and computational linguistics. Word embedding is a text mining technique of establishing relationship between words in textual data (Corpus). The syntactic and semantic meanings of words are realized from the context in which they are used. The concept of distributional hypothesis suggests that words occurring in similar context are semantically similar. Count based embeddings and prediction-based embeddings are the two broad approaches to word embedding. Embeddings capture relationships in language. Embeddings are dense vector representations [9] of the characters.

Embedding layer converts integer indices to dense vectors of length 128.

Input dimension: Size of the vocabulary, which is the number of most frequent words.

Output dimension: Dimension of the dense embedding. It is the vector space in which words will be embedded.

Input length: Length of input sequences which is max length.

Word embeddings are dense vectors with much lower dimensionality. Secondly, the semantic relationships between words are reflected within the distance and direction of the vectors. it's a representation of text where words that have the identical meaning have an analogous representation. In other words, it represents words in an exceedingly system where related words, supported a corpus of relationships, are placed closer together. within the deep learning frameworks like TensorFlow, Keras, this part is typically handled by an embedding layer which stores a lookup table to map the words represented by numeric indexes to their dense vector representations.

D. Embedding Layer

An embedding layer, for lack of a higher name, maybe a word embedding that's learned jointly with a neural network model on a particular linguistic communication processing task, like language modeling or document classification. It requires that document text be cleaned and ready such each word is one-hot encoded. the scale of the vector space is specified as a part of the model, such as 50, 100, or 300 dimensions. The vectors are initialized with small random numbers. The embedding layer is used on the front end of a neural network and is fit in a supervised way using the Backpropagation algorithm. This approach of learning an embedding layer requires a lot of training data and can be slow, but will learn an embedding both targeted to the specific text data and the NLP task.

E. Using Word Embedding

You have some options when it comes time to using word embeddings on your natural language processing project.

- 1. Learn an Embedding
 - You may choose to learn a word embedding for your problem. This will require a large amount of text data to ensure that useful embeddings are learned, such as millions or billions of words. You have two main options when training your word embedding:
 - Learn it Standalone, where a model is trained to be told the embedding, which is saved and used as an element of another model for your task later. this is often a decent approach if you'd prefer to use the identical embedding in multiple models.
 - Learn Jointly, where the embedding is learned as a part of an oversized task-specific model. this is often a decent approach if you simply shall use the embedding on one task.
- 2. Reuse an Embedding
 - It is common for researchers to make pre-trained word embeddings available for free, often under a permissive license so that you can use them on your own academic or commercial projects. For example, both word2vec and Glove word embeddings are available for free download. These are often used on your project rather than training your own embeddings from scratch. you have got two main options when it involves using pre-trained embeddings.
 - Static, where the embedding is kept static and is employed as a component of your model. this is often an acceptable approach if the embedding may be a good suit for your problem and offers good results. Updated, where the pre-trained embedding is employed to seed the model, but the embedding is updated jointly during the training of the model. this might be an honest option if you're looking to induce the foremost out of the model and embedding on your task.

F. SoftMax

SoftMax function calculates the chances distribution of the event over 'n' different events. generally, a way of claiming, this function will calculate the chances of every target class over all possible target classes. Later the calculated probabilities are helpful for determining the target class for the given inputs. the most advantage of using SoftMax is that the output probabilities range. The range will 0 to 1, and also the sum of all the changes is adequate. If the SoftMax function used for the multi-classification model it returns the chances of every class and also the target class will have a high probability. The formula computes the exponential (e-power) of the given input value and also the sum of exponential values of all the values within the inputs. Then the ratio of the exponential of the input value and also the sum of exponential values is that the output of the SoftMax function. it's Used for the multi-classification task and within the different layers of neural networks. The high value will have a better probability than other values. A neural network could also be attempting to work out if there's a dog in a picture. it should be able to produce a probability that a dog is, or is not, within the image, but it might do so individually, for every input. A SoftMax layer, allows the neural network to run a multi-class function. In short, the neural network will now be able to determine the probability that the dog is within the image, in addition, because the probability that additional objects are included in addition.



SoftMax layers are good at determining multi-class probabilities, however, there are limits. SoftMax can become more expensive as the number of classes grows. In those situations, candidate sampling can be a more effective workaround. With candidate sampling, a SoftMax layer will limit the scope of its calculations to a particular set of classes. For example, when determining if an image of a bowl of fruit has apples, the probability does not need to be calculated for every type of fruit, just the apples. Additionally, a SoftMax layer assumes that there is only one member per class, and in situations where an object belongs to multiple classes, a SoftMax layer will not work. In that case, the alternative is to use multiple logistic regressions instead. Properties of SoftMax Function

- The calculated probabilities will be in the range of 0 to 1.
- The sum of all the probabilities is equals to 1.
- Used in multiple classification logistic regression model.
- In building neural networks SoftMax functions used in different layer level.

G. Algorithm

A step-by-step process for how RNN can be implemented using LSTM architecture

- Load in and visualize the data
- Data processing Remove Punctuation
- Tokenize Encode the words and labels
- Training, Validation, Test Dataset Split
- Define the LSTM Network Architecture (Building Model)
- Training the Network
- Testing (on Test data and User-generated data)

H. Working of LSTM Network

- 1. Take input the current input, the previous hidden state and the previous internal cell state.
- 2. Calculate the values of the four different gates by following the below steps: -
 - For each gate, calculate the parameterized vectors for the current input and the previous hidden state by element-wise multiplication with the concerned vector with the respective weights for each gate.
 - Apply the respective activation function for each gate element-wise on the parameterized vectors. Below given is the list of the gates with the activation function to be applied for the gate.
- 3. Calculate the current internal cell state by first calculating the element-wise multiplication vector of the input gate and the input modulation gate, then calculate the element-wise multiplication vector of the forget gate and the previous internal cell state and then adding the two vectors.

$$c_t = i \odot g + f \odot c_{t-1}$$

4. Calculate the current hidden state by first taking the element-wise hyperbolic tangent of the current internal cell state vector and then performing element wise multiplication with the output gate.

$$h_t = o \odot tanh(c_t)$$

Just like Recurrent Neural Networks, an LSTM network also generates an output at each time step and this output is used to train the network using gradient descent.

The only main difference between the Back-Propagation algorithms of Recurrent Neural Networks and Long Short-Term Memory Networks is related to the mathematics of the algorithm.

IV. EXPERIMENT RESULTS

A. Summary of Dataset

For our experimental study we use the IMDB and Amazon Product datasets. IMDB is the large movie review dataset and is a bench mark for movie review dataset that contains a total of 50,000 reviews out of which 25000 are positively polarized and 25000 are negatively polarized. Among the total available reviews, 50,000 reviews are used for training and the remaining 23500 are used for evaluating the performance of the trained model. The objective of this work is to identify the polarity of the given review that is whether the review given is of positive sentiment or negative sentiment. Table 1. Summary of the IMDB dataset

Dataset	Total Samples	Train Samples	Test Samples	Classes
IMDB	50000	25000	25000	2
Amazon	50000	25000	25000	2

B. Layers of our model

Model: "sequential_1"

Layer (type)	Output	Shape	Param #
embedding_1 (Embedding)	(None,	2469, 100)	12129000
lstm_1 (LSTM)	(None,	128)	117248
dense_1 (Dense)	(None,	5)	645
Total params: 12,246,893 Trainable params: 12,246,893 Non-trainable params: 0			

None

C. Model Architecture

We initialize the word embedding layer with random values. Each word is represented with an embedding vector of size 100. The top 6000 words are used in the vocabulary and rare words are removed from the dictionary to avoid unnecessary computations. During training, the hyperparameters that resulted in the best performance are: Dropout is applied with a rate of 0.2. Adam optimizer is used to optimize the model and sparse_categorical_crossentropy is used as the loss function. A batch size of 500 is adopted.

Configuration of the model	Epochs	LSTM Units	Accuracy
	1	128	56.14%
	2	128	75.91%
Embedding	3	128	88.13%
Layer	4	128	93.44%
+	5	128	96.16%
LSTM Layer	6	128	97.67%
+	7	128	98.59%
Dense Layer	8	128	99.11%
	9	128	99.36%
	10	128	99.56%

V. CONCLUSION AND FUTURE SCOPE

In this paper, we have proposed a sentiment classification approach based on LSTM for text data. Users from all over the world express and publicly share their opinions on different topics. Manual analysis of large amounts of such data is very difficult, so a reasonable need for their computer processing has emerged. Sentiment analysis processes people's opinions and attitudes toward products, services, politics, social events, and company strategies. Reviews (from sources such as TripAdvisor, Amazon, and IMDB) and social network posts (mostly from Twitter and Facebook) are categories of textual documents that are the most interesting for sentiment analysis. DL methods such as LSTM show better performance of sentiment classification with 85% accuracy when there are more amounts of training data.

In future we are planning to extend this study to a larger extent where different embedding models can be considered on large variety of the datasets.

REFERENCES

- Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., and Potts, C. (2011, June). Learning word vectors for sentiment analysis. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1 (pp. 142-150). Association for Computational Linguistics
- [2] ApoorvAgarwal, BoyiXie, Ilia Vovsha, Owen Rambow, Rebecca Passonneau. Sentiment Analysis on twitter data. Department of Computer Science, Columbia University, New York, NY 10027 USA.
- [3] https://doi.org/10.1007/s12088-011-0245-8
- Hoch Reiter S, Schmidhuber, J. (1997). Long short-term memory. In: Neural Computation 9(8): 1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735
- [5] PRANJAL SRIVASTAVA, Essentials of Deep Learning: Introduction to Long Short Term Memory, DECEMBER 10, 2017.
- [6] http://ai.stanford.edu/~amaas/data/sentiment/
- Ain QT, Ali M, Riazy A, Noureenz A, Kamranz M, Hayat B, Rehman A. (2017). Sentiment analysis using deep learning techniques: A review. In: International Journal of Advanced omputer Science and Applications (IJACSA) 8(6): https://doi.org/10.14569/IJACSA.2017.080657
- [8] Sokolova M. (2018). Big text advantages and challenges: Classification perspective. In: International Journal of Data Science and Analytics 5(1): 1-10. https://doi.org/10.1007/s41060-017-0087-5
- [9] https://towardsdatascience.com/understanding-confusionmatrixa9ad42dcfd62
- [10] Machine Learning Tom M. Mitchell McGraw-Hill