

Machine Learning Engineer Nanodegree

Capstone Project

Ranveer Dutta

September 8th, 2018

I. Definition

Project Overview

Sentiment analysis is contextual mining of text which identifies and extracts subjective information in source material, and helping a business to understand the social sentiment of their brand, product or service while monitoring online conversations.

The motive of this project is to understand movie rating based on the review people have provided in texts.

Most sentiment prediction systems work just by looking at words in isolation, giving positive points for positive words and negative points for negative words and then summing up these points. That way, the order of words is ignored and important information is lost.

Let's look at the below example statement:

"This movie was actually neither that funny, nor super witty"

A naive sentiment analysis tool that predicts based on individual words may consider the above statement as positive since it has 2 positive words "**funny**" and "**witty**" but actually the entire statement is negative. In contrast, our new deep learning model actually builds up a representation of whole sentences based on the sentence structure. It computes the sentiment based on how words compose the meaning of longer phrases. This way, the model is not as easily fooled as naive models. For example, our model learned that funny and witty are positive but the complete sentence is still negative overall.

The expected output is to label phrases on a scale of five values: negative, somewhat negative, neutral, somewhat positive, positive. Obstacles like sentence negation, sarcasm, terseness, language ambiguity, and many others make this task very challenging.

Problem Statement

The goal of this project is to:

- Create a text parser which will read through the ratings provided for movies in texts.
- Create a model which can infer the meaning of different text statements and outputs sentiments of the reviewer.
- The output data can be used by different clients to define the movie ratings given the text reviews.

Metrics

We have used mean square error (MSE) to observe the performance of the simple CNN estimator during training. At each step of the training, we will calculate MSE.

At the end we'll calculate the accuracy on test set using the trained model.

II. Analysis

Data Exploration

In this study we'll use the Rotten Tomatoes dataset. The dataset is comprised of tab-separated files with phrases. The train/test split has been preserved for the purposes of benchmarking, but the sentences have been shuffled from their original order. Each Sentence has been parsed into many phrases by the Stanford parser. Each phrase has a Phraseld. Each sentence has a Sentenceld. Phrases that are repeated (such as short/common words) are only included once in the data.

- **train.tsv** contains the phrases and their associated sentiment labels. We have additionally provided a Sentenceld so that we can track which phrases belong to a single sentence.
- **test.tsv** contains just phrases. You must assign a sentiment label to each phrase.

Few examples from the training set is as below:

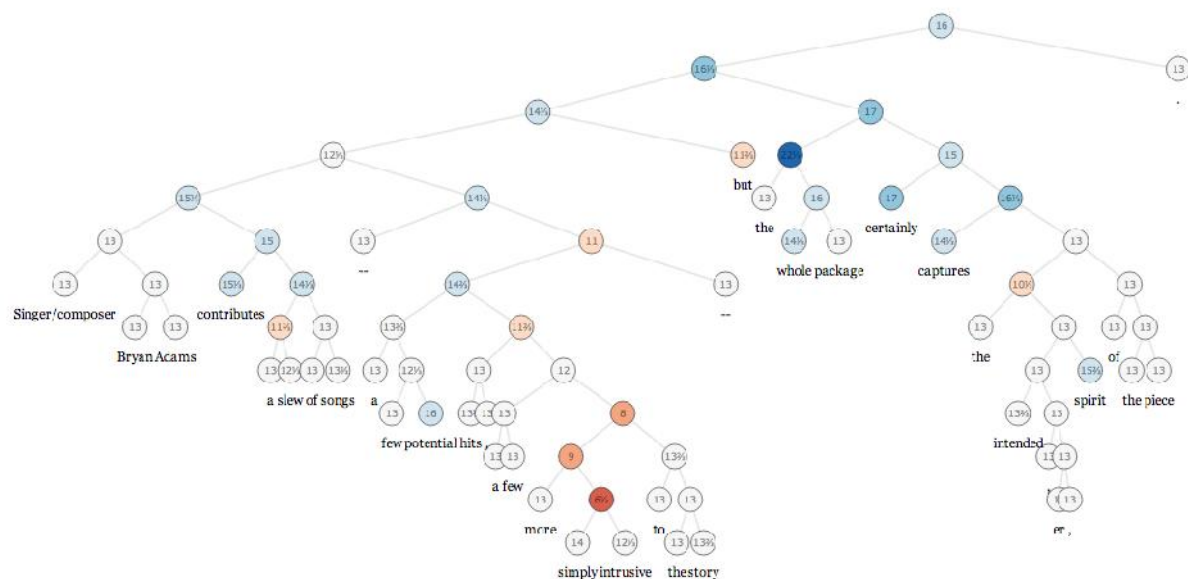
Phraseld	Sentenceld	Phrase	Sentiment (label)
16	1	that what is good for the goose	2
17	1	that	2
22	1	good for the goose	3
34	1	the gander , some of which occasionally amuses but none of which amounts to much of a story	1

In the above example we can see how the same sentence **“A series of escapades demonstrating the adage that what is good for the goose is also good for the gander , some of which occasionally amuses but none of which amounts to much of a story”** have different sentiment when used in parts.

The sentiment labels are:

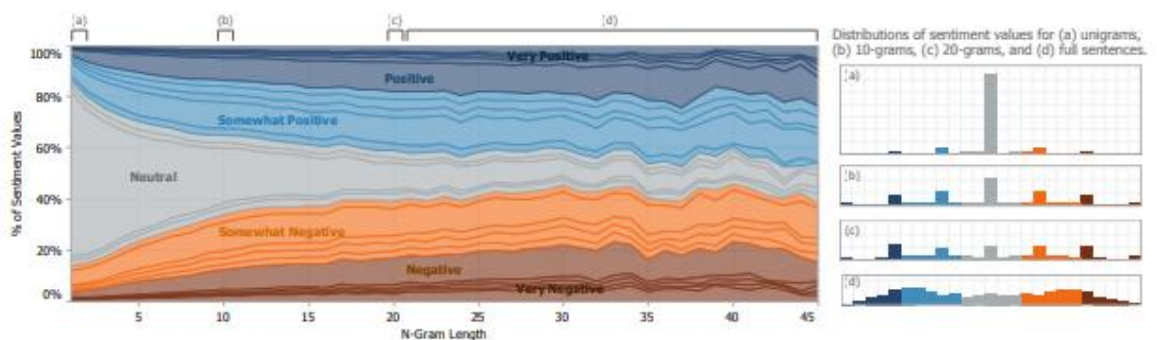
- 0 - negative
- 1 - somewhat negative
- 2 - neutral
- 3 - somewhat positive
- 4 - positive

Below is one example of treebank:



Exploratory Visualization

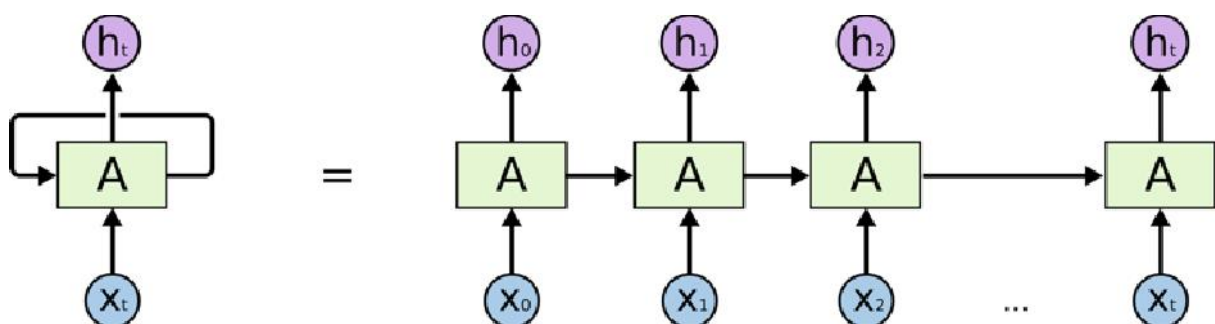
The original dataset includes 10,662 sentences, half of which were considered positive and the other half negative. Each label is extracted from a longer movie review and reflects the writer's overall intention for this review.



Above figure shows the normalized label distributions at each n-gram length. Starting at length 20, the majority are full sentences. One of the findings from labeling sentences based on reader's perception is that many of them could be considered neutral. We also notice that stronger sentiment often builds up in longer phrases and the majority of the shorter phrases are neutral. Another observation is that most annotators moved the slider to one of the five positions: negative, somewhat negative, neutral, positive or somewhat positive. The extreme values were rarely used and the slider was not often left in between the ticks. Hence, even a 5-class classification into these categories captures the main variability of the labels

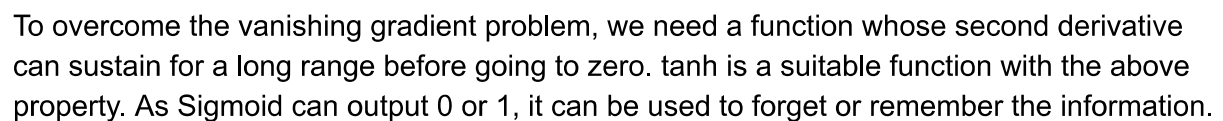
Algorithms and Techniques

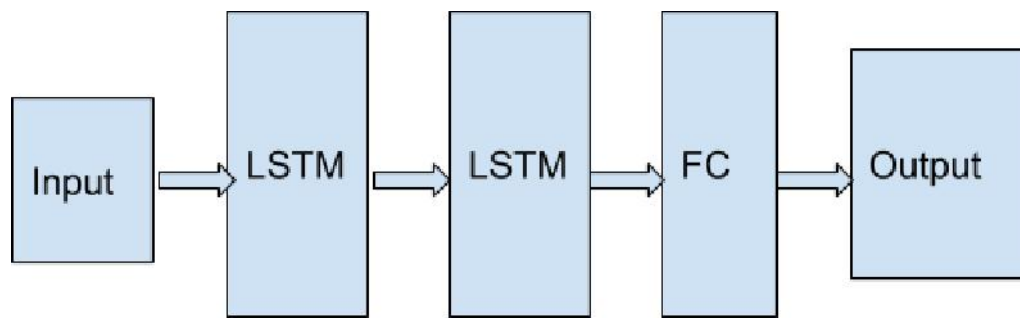
As stated in the problem statement, we need to extract the meaning of the entire statement. For that purpose, our neural network must consider events from previous steps. For this purpose Recursive Neural Network is good choice.



During the training of RNN, as the information goes in loop again and again which results in very large updates to neural network model weights. This is due to the accumulation of error gradients during an update and hence, results in an unstable network. At an extreme, the values of weights can become so large as to overflow and result in NaN values. The explosion occurs through exponential growth by repeatedly multiplying gradients through the network layers that have values larger than 1 or vanishing occurs if the values are less than 1.

The diagram illustrates the internal structure of an LSTM cell. It shows the flow of information through four gates (sigmoid and tanh) and the cell state. The inputs are $c(t-1)$ and $h(t-1)$, and the current input x_t is also shown. The cell state $c(t)$ is updated by adding the product of the forget gate and the previous cell state to the product of the input gate and the candidate cell state. The hidden state h_t is then calculated by applying a tanh activation to the cell state and multiplying it by the output gate.





Benchmark

To Benchmark the solution, I kept few dataset of Movie review text from we website www.bookmyshow.com. It has movie review as well. In the implementation the provided dataset is already parsed and presented in the required format. For benchmarking, I had manually parsed and created the dataset. My expectation was to get 80-85% accuracy.

III. Methodology

Data Preprocessing

The data was provided in 3 separate files.

File name	Columns	Preprocessing done
train.tsv	Phraseld, Sentenceld, Phrase, Sentiment	<ol style="list-style-type: none"> 1. Phraseld and Sentenceld is dropped 2. All the characters in Phrase are converted into lowercase 3. Text from Phrase are converted to sequences using <code>keras tokenizer.texts_to_sequences</code> 4. The sentiment indexes are converted to vector using <code>to_categorical</code> Of keras. 5. Divided the train data into training set(80%) and validation set(20%).

test.tsv	Phraseld, Sentenceld, Phrase	<ol style="list-style-type: none"> 1. Phraseld and Sentenceld is dropped 2. All the characters in Phrase are converted into lowercase 3. Text from Phrase are converted to sequences using <code>keras tokenizer.texts_to_sequences</code>
sampleSubmission.csv	Phraseld, Sentiment	<ol style="list-style-type: none"> 1. Phraseld is dropped 2. The sentiment indexes are converted to vector using <code>to_categorical</code> Of keras.

Implementation

- The LSTM neural network which is a variant of Recurring Neural Network(RNN) is used.
- We used Keras implementation library for this.
- We used 2 layers of LSTM with varying units are used.
- AT the end one fully connected Dense layer is used which outputs 5 different classes of sentiments.
- To avoid overfitting, dropouts are used.
- At the Dense layer, activation function"softmax" is used.

Refinement

To refine the accuracy, I played around with below parameters:

- 1) Epochs: I started with 5 epochs and at the end settled with 10 epochs.
- 2) Units in LSTM layer: I tried different combination and at the end given 18 and 64.
- 3) I trained the model with 128 batch size.
- 4) The activation function I used is "softmax" at the dense layer.
- 5) I used loss function "categorical_crossentropy"

IV. Results

Model Evaluation and Validation

After the training I got the best accuracy of 74.62%. From 5th epoch onwards the loss value did not improve. When I tested with the test data, I got the accuracy of 59.16%. This is average accuracy and does not look like great performance to be used to real application.

Justification

During benchmarking I got the correct result of close to 65% of data which is not very good. The final solution is not upto the expectation and needs more refinement.

V. Conclusion

Free-Form Visualization

Sort By Popularity

Reviewer	Review Title	Rating	Date	Thumbs Up	Thumbs Down
Vishnu	class movie	★★★★★	15/8/18	(399)	(182)
Sukumar	GOLD	★★★★★	15/8/18	(240)	(48)
Smeet	Inspiring Movie	★★★★★	15/8/18	(153)	(37)

Vishnu
drama, imotion, very nice roll played by akshay and team. nice direction by rima Kagti. it is one gift for nation on independence . that kind of movies boost for our national game hockey

Sukumar
Movie Bahut achi hai jake dekh lo. Akshay sir ne bahut aachi acting ki hai.Maine 1day 1show dekha Din ki shurvut bahu achi hue.Bahut BADIYA👍👍 JAY HINDINININ

Smeet
Fantastic Role of Akshay Kumar.. Owsme.. Heart touching movie.. A great tribute to hockey players. Story and Script was too good.. Must watch this patriotic movie.. :)

The above screen shot is from a only movie ticket booking site. If trained properly, it can easily say that 1st and 3rd belong to class 4. But, the 2nd review is confusing, since the review was written using English script but the semantics is from an Indian language Hindi.

Reflection

The interesting part of the problem was to understand the sentiments by looking at the entire statement than just looking for occurrence of few words. One of complex part is to create the dataset using different phrases of the statement. In this particular problem the parsed dataset was already provided, so after some preprocessing it was ready to work with the solution. Another problem would be to parse reviews which has combination of phrases which belong to different sentiments. Some movies may have first half as entertaining and second half half may be boring. The reviewer may not provide final conclusion and give mixed review in both positive and negative phrases.

Another problem is with semantics. If we travel from one geography to another, then people have different expectation from a good/bad movie and their language of expression may also be different. So having more refined dataset will help that will be mapped to that geography user type.

The final solution is having accuracy of 60% only , but I strongly feel it can improved with similar architecture using more combination of parameters.

Improvement

The solution can be improved further to reach acceptable accuracy level.

- We can try adding different combination of LSTM and Dense layer..
- We can try bidirectional LSTM.
- We can change epoch and see the improvement.
- The size of LSTM can be varies.
- The batch size can be changed.
- We can try different dropout factor and activation function.
- We can try using pooling layer in-between.

References:

<https://www.kaggle.com/c/movie-review-sentiment-analysis-kernels-only>

<https://nlp.stanford.edu/sentiment/>

<https://towardsdatascience.com/understanding-lstm-and-its-quick-implementation-in-keras-for-sentiment-analysis-af410fd85b47>