# Sentiment Analysis of Online Product Reviews

Venkatesh Gopinath, Monica Ram, Mahitha Guduri
*Dept. of Mechanical and Industrial Engineering*
*College of Engineering, Northeastern University*
Boston, MA
bogem.v@northeastern.edu, ram.m@northeastern.edu, guduri.m@northeastern.edu

*Abstract*—Sentiment analysis is a method to extract opinions and attitudes from text data using natural language processing (NLP). The aim of this research project is to perform sentiment analysis on Amazon customer reviews using two different techniques from NLTK: VADER (Valence Aware Dictionary and Sentiment Reasoner) and Roberta Pretrained Model from Hugging Face and three Neural Network models. We will preprocess the text data, perform tokenization, and train the model using the Keras deep learning library.

*Index Terms*—Reviews, VADER, Roberta, Sequential, LSTM, glove embedding, Pad Sequences, Keras.

## I. INTRODUCTION

Sentiment analysis has become increasingly popular with the increasing amount of user-generated content on the internet. This technique is used to extract subjective information from text data sources such as reviews, social media posts, and news articles. Sentiment analysis can be used to understand the overall sentiment of a particular brand or product, analyze customer feedback, and identify trends in customer behavior. In recent years, deep learning techniques, such as neural networks, have shown impressive results in sentiment analysis tasks.

In this paper, we performed few basic exploratory data analysis to understand the distribution of the data. We later performed sentiment analysis using few models from NLTK and neural network models from keras. The study uses the Amazon review dataset for product reviews. We utilized two approaches from NLTK to perform sentiment analysis: VADER (Valence Aware Dictionary and Sentiment Reasoner) and Roberta Pretrained Models. The neural network model used in this project is a sequential model from Keras with an embedding layer, an LSTM layer, and a dense output layer. We conducted a study to compare the performance of these methods.

We will use the Amazon customer reviews dataset to train and test our model. The Amazon reviews dataset is a collection consisting of 3,744,047 reviews on a variety of products. It contains product reviews as text and has a score which is star reviews from 1 to 5 with 5 being most liked and 1 being least liked.

## II. BACKGROUNDS

With the advances in Deep Learning in parallel with higher computation capabilities using cloud computing techniques, approaches like Recurrent Neural Networks, Long Short-Term Memory and the latest Transformer based methods have been a huge success. One such approach of text-based sentiment analysis, Recurrent Neural Networks is considered the most powerful method.

## III. APPROACH

### A. Dataset

The Amazon reviews dataset [1] consists of features like 'Id', 'Product Id', 'User Id', 'Profile Name', 'Helpfulness Numerator', 'Helpfulness Denominator', 'Score', 'Time', 'Summary', 'Text'. We will be working on 'Text', 'Id', and 'Score'.The dataset is split into training, validation, and test datasets. The training dataset contains 60% of the data, the validation dataset contains 20% of the data and the test dataset contains the remaining 20% of the data. Then the test dataset is further split into validation andtesting sets with stratified random sampling to ensure the classes are balanced in all the training, validation, and test sets.

### B. Pre-processing

- In Exploratory Data Analysis (EDA), for each label in the train set to get an overviewof the most frequent words. Many HTML tags like" p"," href", etc. are identified. These words are removed.

- We later plotted a bar plot to compare the count of reviews based on the star reviews they received. We observed that most of the reviews have a five-star rating. And the count decreased gradually with the ratings.

- We wrote a function to perform a few sorts of pre-processing like:
  - Lowercase the entire text
  - remove all html tags like "p", "href", etc. from dataset
  - Replace special characters and punctuation with a blank (except for full stop(.))
  - split the text into words
  - remove stop words like to, an, the, of, over, his, her…

- Punctuation, numbers are removed and converted into tokens with" " as a delimiter.
- In the same function we also collected pos_tags for all the words and lemmatized them using WordNetLemmatizer ().
- We also plotted a bar plot to compare the count of reviews and the length of the words in them.
- We observed that most of the reviews have less than 200 words. Hence, we put a threshold of 200 words and considered the reviews only whose length is less than 200.

## C. Vector Embeddings
- We performed embedding of the text using pre-trained word vectors, 'glove.6B.100d.txt' which is publicly available at the Stanford website.
- The dimension of each word embedding is set to 100.
- We then created an embedding matrix for directly feeding to Embedding () layer of Keras.
- We wrote a function to average out embedding for each corresponding word in a review.

## D. Loss Function
Since our problem is Binary (Positive and Negative), we will be using Binary Cross Entropy Loss is used.

$$L_{binarycrossentropy}(y, p) = \text{abs} (Y\_pred - Y\_actual)$$

## E. Evaluation Metrics
We chose to use Accuracy as a performance measure of our Train, Validation, and Test sets.



Fig.1. Evaluation metrics VADER Polarity Scores



Fig.2. Evaluation metrics Roberta Polarity Scores



Fig.3. Evaluation for average embedding



Fig. 4. Evaluation for average glove embedding

## IV. MODELS

### A. VADER Model from NLTK library
We used NLTK's Sentiment Intensity Analyzer to get the neg/neu/pos scores of the text. This uses a "bag of words" approach and stop words are removed. Also, each word is scored and combined to a total score.



Fig. 5. VADAR Model

### B. Roberta Pretrained Model from Hugging face library
It builds on BERT and modifies key hyperparameters, removing the next-sentence pretraining objective and training with much larger mini-batches and learning rates.



Fig. 6. Roberta Model

## C. Sequential model from Keras

The sequential model allows us to create layer-by-layer for most problems. It is limited. It does not allow us to create models that share layers or have multiple inputs or outputs.

The same model is used for both average word embedded and average glove-embedded data frames.

```
1  model_avg_embedding = Sequential()
2
3  model_avg_embedding.add(Dense(250, input_shape=(EMBEDDING_DIM, ), activation='relu'))
4  model_avg_embedding.add(Dropout(0.3))
5
6  model_avg_embedding.add(Dense(100, activation='relu'))
7  model_avg_embedding.add(Dropout(0.4))
8
9  model_avg_embedding.add(Dense(1, activation='sigmoid'))
10
11 optimizer = Adam(lr=1e-3, beta_1=0.9, beta_2=0.999, epsilon=1e-08, decay=0.0, amsgrad=False)
12
13 model_avg_embedding.compile(loss='binary_crossentropy', optimizer=optimizer, metrics=['accuracy'])
14 model_avg_embedding.summary()
Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
dense (Dense)                (None, 250)               25250
_____
dropout (Dropout)            (None, 250)               0
_____
dense_1 (Dense)              (None, 100)               25100
_____
dropout_1 (Dropout)          (None, 100)               0
_____
dense_2 (Dense)              (None, 1)                 101
=================================================================
Total params: 50,451
Trainable params: 50,451
Non-trainable params: 0
```

Fig. 7. Sequential Model for Word Embedded

```
1
2  model_pretrained_glove = Sequential()
3
4  embedding_layer = Embedding(len(tokenizer.word_index)+1,
5                              EMBEDDING_DIM,
6                              embeddings_initializer=Constant(embedding_matrix),
7                              input_length=MAX_SEQUENCE_LENGTH,
8                              trainable=False)
9  model_pretrained_glove.add(embedding_layer)
10 model_pretrained_glove.add(Dropout(0.3))
11
12 model_pretrained_glove.add(Conv1D(64, 5, activation='relu'))
13 model_pretrained_glove.add(MaxPooling1D(pool_size=4))
14
15 model_pretrained_glove.add(LSTM(100))
16 model_pretrained_glove.add(Dropout(0.4))
17
18 model_pretrained_glove.add(Dense(1, activation='sigmoid'))
19
20 optimizer = Adam(lr=1e-3, beta_1=0.9, beta_2=0.999, epsilon=1e-08, decay=0.0, amsgrad=False)
21
22 model_pretrained_glove.compile(loss='binary_crossentropy', optimizer=optimizer, metrics=['accuracy'])
23 model_pretrained_glove.summary()
Model: "sequential_1"
_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding (Embedding)        (None, 200, 100)          271100
_____
dropout_2 (Dropout)          (None, 200, 100)          0
_____
conv1d (Conv1D)              (None, 196, 64)           32064
_____
max_pooling1d (MaxPooling1D) (None, 49, 64)            0
_____
lstm (LSTM)                  (None, 100)               66000
_____
dropout_3 (Dropout)          (None, 100)               0
_____
dense_3 (Dense)              (None, 1)                 101
=================================================================
Total params: 369,265
Trainable params: 98,165
Non-trainable params: 271,100
```

Fig. 8. Sequential Model for glove Embedded

## V. RESULTS

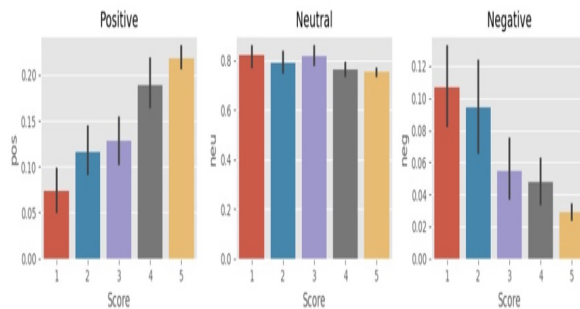### A. VADAR Results



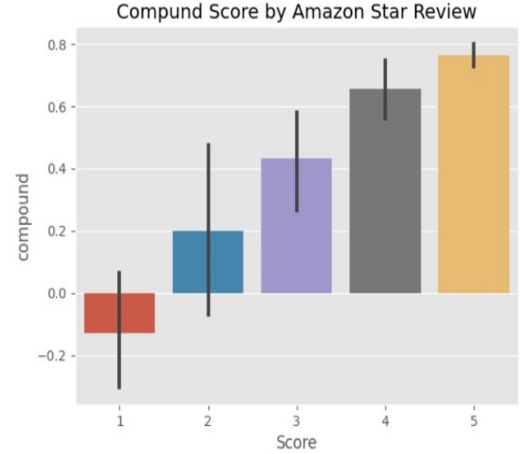Fig. 9. VADAR Score Vs Star Rating



Fig. 10. VADAR Compound Score VS Star Rating
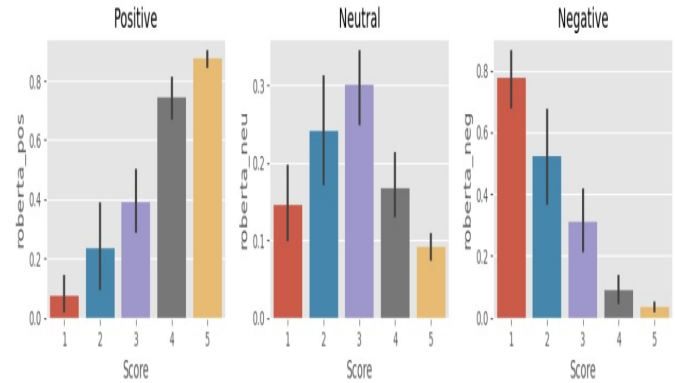
### B. Roberta Results



Fig. 11. Roberta Score Vs Star Rating

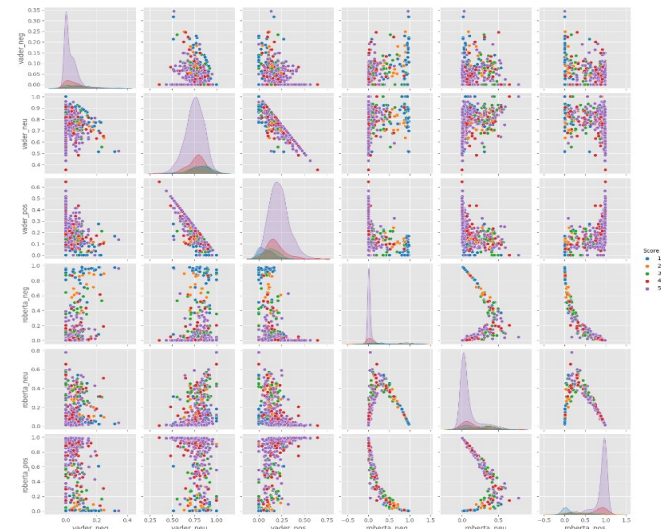### C. Comparing Roberta and VADER Results



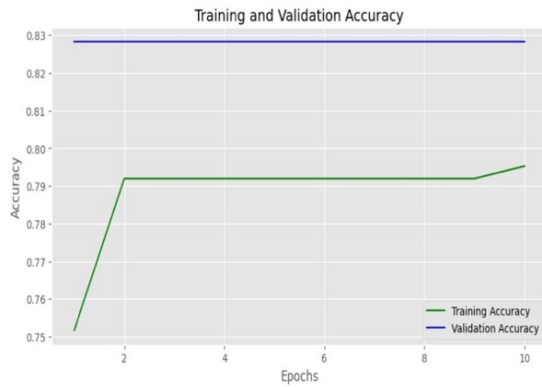Fig. 12. Roberta Vs VADER Results

## D. Sequential results


Fig. 13. Accuracy Vs Epochs for Train & Validation

```
1 # evaluating the model on test dataset
2 model_avg_embedding.evaluate(x_avg_test, y_test)
```

4/4 [==============================] - 0s 2ms/step - loss: 0.3705 - accuracy: 0.8800

[0.37054443359375, 0.8799999952316284]

Fig. 14. Loss & Accuracy Score for Word Embedded


Fig. 15. Loss Vs Epochs for Glove Embedded

```
1 # evaluating the model on test dataset
2 model_pretrained_glove.evaluate(x_test, y_test)
```

4/4 [==============================] - 0s 15ms/step - loss: 0.4047 - accuracy: 0.880

[0.4046874940395355, 0.8799999952316284]

Fig. 16. Loss & Accuracy Scores for Glove Embedded

## VI. CONCLUSION

Based on our comparison between VADER and Roberta Models, we can tell that Roberta has performed a lot better.

Also, between word embedded Sequential model and Glove embedded Sequential model, later resulted in better scores

The next steps could be

- To use a Recurrent neural network which works best for text-based models.

- We can extend our NLTK model comparison to transform that token comparison to transformers that tokenises automatically.

REFERENCES

[1] https://keras.io/api/models/sequential/
[2] https://www.nltk.org/api/nltk.sentiment.vader.html#module-nltk.sentiment.vader

[3] Dataset: https://www.kaggle.com/datasets/snap/amazon-fine-food-reviews