# SENTIMENT ANALYSIS ON AMAZON REVIEWS.

- Raj Thaker.

rpthaker@usc.edu

CSCI 561: FOUNDATIONS OF ARTIFICIAL INTELLIGENCE.
HOMEWORK – 2

DATE : 07/26/2018

➢ <u>Dataset</u>

- The given dataset is an extract from the Amazon Reviews Kaggle competition.
- The dataset of 400,000 reviews is in the form of CSV file delimited by '|'.
- Data contains 2 fields – Classification of Review - (Positive/Negative) and corresponding Review.

➢ <u>Goal</u>

- The goal is to perform sentiment analysis to determine whether a review is positive or negative.
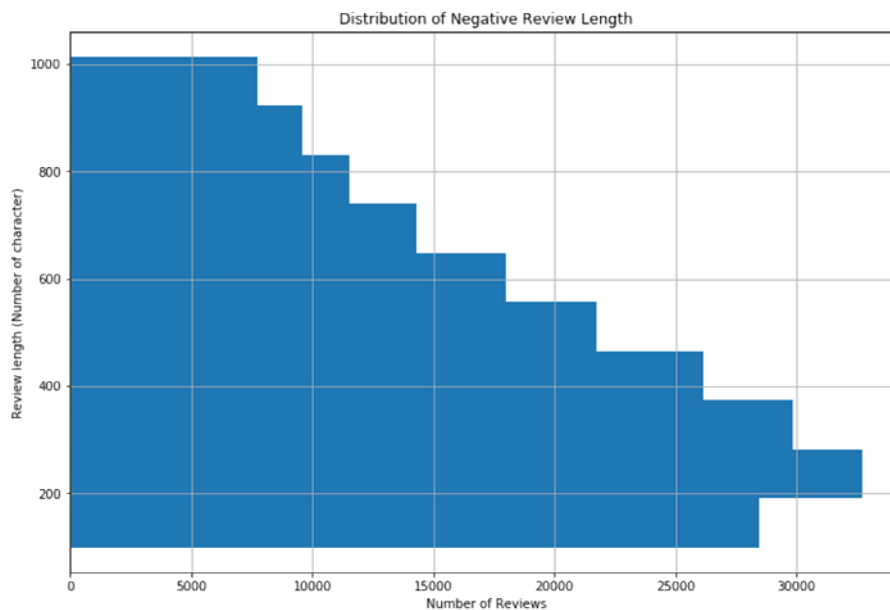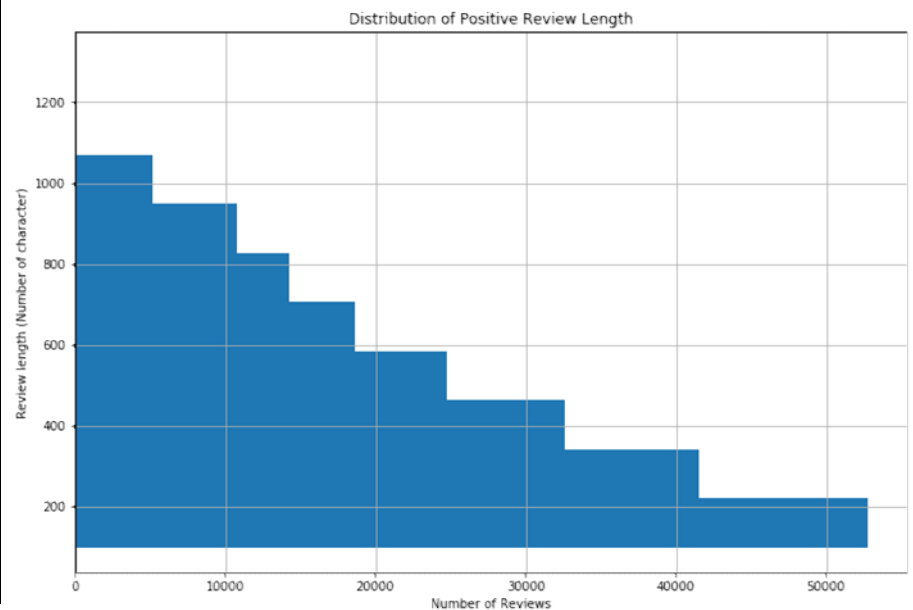
Reference for Kaggle Competition - https://www.kaggle.com/bittlingmayer/amazonreviews

➢ <u>Exploratory Data Analysis.</u>

• The given data contains ***equal number*** of positive and negative reviews.

• Reviews Length Distribution:

Length Distribution of Negative Reviews. | Length Distribution of Positive Reviews.



Distribution of Negative Review Length



Distribution of Positive Review Length

➤ Exploratory Data Analysis.

Positive Words Cloud.

Negative Words Cloud

• Most reviews spoken about:

➢ <u>Data Preprocessing.</u>

- Removal of Punctuations, Stop words, Stemming.
- Count Vectorizer.

➢ <u>How data was setup for training the model.</u>

- To understand different ML Classification algorithms, the given dataset was divided in following ways –
    - I.      50K records were trained | 80K records tested.
    - II.     100K records were trained | 80K records tested.
    - III.    200K records were trained | 80K records tested.
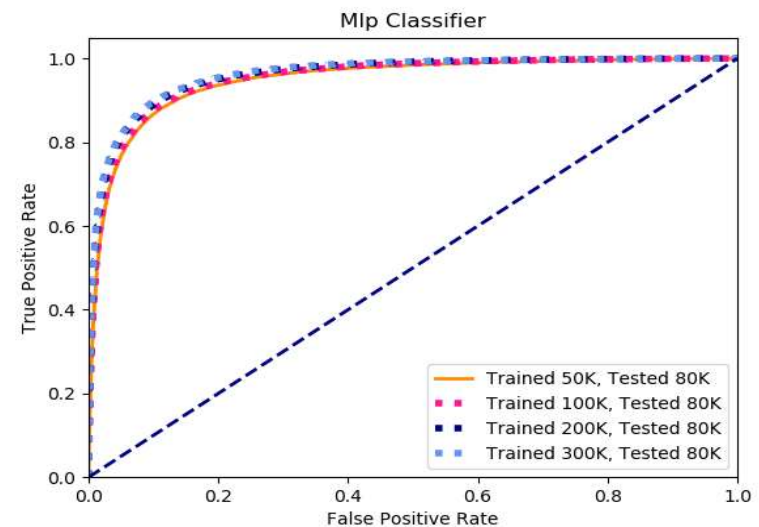    - IV.    310K records were trained | 80K records tested.

➢ <u>3 Machine Learning Algorithms Used:</u>

1. Neural Network – MLPClassifier.
2. Decision Tree
3. Logistic Regression

1. Neural Network – MLPClassifier.

|  | 50K Train \| 80K Test | 100K Train \| 80K Test | 200K Train \| 80K Test | 310K Train \| 80K Test |
|---|---|---|---|---|
| Accuracy | 88.40 | 88.85 | 89.49 | 89.85 |
| Precision | 0.87 | 0.88 | 0.89 | 0.89 |
| Recall | 0.89 | 0.89 | 0.88 | 0.90 |



Mlp Classifier

True Positive Rate vs False Positive Rate

Trained 50K, Tested 80K
Trained 100K, Tested 80K
Trained 200K, Tested 80K
Trained 300K, Tested 80K

> ## 3 Machine Learning Algorithms Used:

## 2.  Decision Tree

|  | 50K Train \| 80K Test | 100K Train \| 80K Test | 200K Train \| 80K Test | 310K Train \| 80K Test |
|---|---|---|---|---|
| Accuracy | 75.41 | 76.33 | 76.51 | 77.19 |
| Precision | 0.75 | 0.76 | 0.76 | 0.77 |
| Recall | 0.75 | 0.76 | 0.76 | 0.76 |



Decision Tree.

## 3.  Logistic Regression

|  | 50K Train \| 80K Test | 100K Train \| 80K Test | 200K Train \| 80K Test | 310K Train \| 80K Test |
|---|---|---|---|---|
| Accuracy | 87.40 | 88.52 | 89.23 | 77.19 |
| Precision | 0.87 | 0.88 | 0.89 | 0.89 |
| Recall | 0.87 | 0.88 | 0.88 | 0.89 |



Logistic Classifier

➢ <u>Final Verdict – Neural Networks!!</u>



ROC Curves for 3 Classifiers.

Neural Network. AUC : 0.898583519822
Decision Tree.  AUC : 0.7719703463968152
Logistic Regression.  AUC : 0.8953389714479478