# Project Report

# Sentiment Analysis on Amazon Fine Food Reviews

**Version 1.0 approved**

**Prepared by Aakash Rami**

**Pratik Mulye**

**Shubham Heda**

**Tanisha Mandvikar**

**Stevens Institute of Technology**

**25 November 2019**

**Sentiment Analysis on Amazon Fine Food Reviews**

**Table of Content**

**Sentiment Analysis on Amazon Fine Food Reviews**

# Business Understanding

Amazon is a multinational e-commerce company that generates $207 billion from retail operations in the United States and abroad. Amazon dominates the retail e-commerce business with about 50% of all online sales in the U.S. in 2018 as per eMarketer. Amazon e-commerce has also launched an Amazon Fine Foods services with the purpose of delivering food services to their customers.

As professor discussed in one of his class, the way 'Target' used data analysis to predict a girl was pregnant motivated us to practically get consumer insights for some organization and help them implement customer targeted marketing. For our project we chose a dataset on Amazon Fine Foods review which has more than 500,000 customer reviews. As this dataset is huge and provides feedback in form of score i.e. on scale of one to five & comments in form of text, we will be trying to generate customer behavior & insights using sentimental analysis.

Amazon Fine Foods consists of half a million reviews and our analysis includes:

- Identifying customer comments and predicting helpfulness of the reviews according to score
- Classifying consumer comments into positive or negative sentiment
- Build word cloud to analyze customer reviews based on product score
- Build predictor models based on customer comments

**Sentiment Analysis on Amazon Fine Food Reviews**

This analysis approach will help us gain insights about customer behavior and help us into future deployment of customer targeted marketing. It will also allow the company to only deliver customer-oriented food products which will prove to be cost efficient.

## Data Analytic Solution:

After loading the dataset, we selected columns that are useful towards our analysis and we analyzed the reviews. The reviews are skewed to positive or negative, more than half of the reviews are with '0' votes for usefulness, also many people agree with score 5 reviews as shown in heatmap. We will then need to pre-process this data. Then we will categorize them into positive, neutral and negative. Based on such analysis it will help the company to focus on specific products. It brings a clear picture to the organization where we find which food products customer liked the most and which food products they disliked. Accordingly, they will focus on production of customer-oriented product only rather than producing every other food product which as result will minimize food waste, provide a cost-effective measure and more profit to Amazon.

After exploring the reviews based on score and word cloud it was necessary to dive deep into text mining on text reviews. So, we developed a text-based classification system that can accurately predict the positivity of Amazon online consumer reviews. Text mining was used to perform a binary classification using the combination of text-based features and machine learning classification algorithms. The binary classes will be predicted as - '1' being 'Positive' and '0' being 'Negative' in column named 'Positivity'. Text-based features will be used for this analysis which includes features extracted from review text. The prediction model will be obtained based on the available training data which consists of the customer review text, summary, score, details

**Sentiment Analysis on Amazon Fine Food Reviews**

on the helpfulness votes and more information. The algorithms that will be used for classification purpose are- Multinomial Logistic Regression, Decision Tree, Random Forest. Text pre-processing is also needed to avoid overfitting and reduce computational complexity by removing common stop words using 'nltk' library.

*(Source: Text Analytics for Beginners using NLTK)*

Also, the helpfulness measure is determined based on the number of users who voted the review as 'helpful'.

## Business Value:

The target value of our project is to find out the most reviewed product and then we categorize them in positive and negative. This categorization can help us get deep insights about the products which are not performing satisfactorily, and the company can focus particularly on these products. They can decide on the products which need to be manufactured or the ones whose production should be ceased at once and thus enhancing inventory management. This will help in increasing revenues and decrease the cost required for production & stocking of items which have low income. Also, sentiment will focus on number of user as well. Based on customers rating and reviews Amazon will recommend similar food services.

**Sentiment Analysis on Amazon Fine Food Reviews**

# Data Understanding:

The dataset is acquired from SNAP-Stanford university. It consists of half a million reviews for fine foods from Amazon and other products from different categories as well. Reviews include product and user information, ratings and a text review. The 10 columns in the dataset are as follows:

- Id
- UserId
- ProductId
- ProfileName
- HelpfulnessNumerator
- HelpfulnessDenominator
- Score
- Time
- Summary
- Reviews

| Data Statistics | |
|---|---|
| Number of reviewers | 568,454 |
| Number of users | 256,059 |
| Number of products | 74,258 |
| Users with > 50 reviews | 260 |
| Median number of words per review | 56 |
| Time Span | Oct 1999 – Oct 2012 |

*(Data Source: SNAP: Web data: Amazon Fine Foods reviews. (2019). Snap.stanford.edu.)*

**Sentiment Analysis on Amazon Fine Food Reviews**

## Data Preparation:

Our focus is to analyze the Amazon's customer reviews, identify the positive or negative reviews, based on that we will perform analysis to improve the business's product and increase the number of customers.

To achieve this first we will download the dataset. In which we have columns like, Product ID, Username, score, reviews etc.

Among these features, we removed unnecessary columns (which will not be useful while developing the model). To achieve this, we performed 'pearson' correlation among all the features and we analyzed that correlation of feature 'Time' with respect to 'Id' is less so we dropped that column/feature.

|  | Id | HelpfulnessNumerator | HelpfulnessDenominator | Score | Time |
|---|---|---|---|---|---|
| Id | 1.000000 | 0.001227 | 0.000770 | 0.010706 | 0.007912 |
| HelpfulnessNumerator | 0.001227 | 1.000000 | 0.974689 | -0.032590 | -0.154818 |
| HelpfulnessDenominator | 0.000770 | 0.974689 | 1.000000 | -0.097986 | -0.173289 |
| Score | 0.010706 | -0.032590 | -0.097986 | 1.000000 | -0.062760 |
| Time | 0.007912 | -0.154818 | -0.173289 | -0.062760 | 1.000000 |

**Fig. 3. (a)**

Then we checked for null values and removed them.

```
Id                        0
ProductId                 0
UserId                    0
ProfileName              16
HelpfulnessNumerator      0
HelpfulnessDenominator    0
Score                     0
Time                      0
Summary                  27
Text                      0
dtype: int64
```

**Fig. 3. (b) Before removing null values**

**Sentiment Analysis on Amazon Fine Food Reviews**

```
Id                       0
ProductId                0
UserId                   0
HelpfulnessNumerator     0
HelpfulnessDenominator   0
Score                    0
Summary                  0
Text                     0
dtype: int64
```

**Fig 3. (c) After removing null values**



```
In [15]: import seaborn as sns
         sns.pairplot(df)
Out[15]: <seaborn.axisgrid.PairGrid at 0x299a91a7160>
```

```
In [17]: df.shape
Out[17]: (568454, 10)
```

**Fig. 3. (d)**

**Sentiment Analysis on Amazon Fine Food Reviews**

Then we categorized the reviews provided by customers into three categories: Positive, neutral and negative based on the score given by the customers with respect to that review. If the score is more than 3, then it is considered as positive. If it is 3 then its neutral and if its below 3 then its negative i.e. customer is not satisfied with food.

From below graph we can easily analyze the number of positive, negative and neutral feedbacks. Based on following graph, we can see that customers has given majority positive feedback, which seems they are happy with the services and customers also liked the food.

```
In [14]: sns.countplot(foods1_na['Sentiment'])
         plt.show()
```



**Fig. 3. (e)**
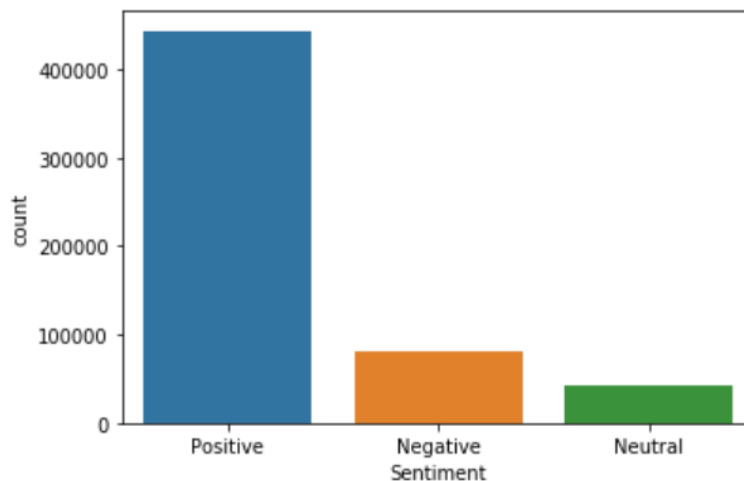
We also had generated HeatMap to identify which category has received most reviews. Heatmap is another way to analyze reviews. In Heatmap we can divided it into percentage and analyzed how much percentage of reviewers has given score 1 or 2 or 3 and so on and how many consumers find it useful. For reference please see below screenshot.

**Sentiment Analysis on Amazon Fine Food Reviews**

From the figure 3.(f) we can see that reviews are skewed towards positive. More than half of the reviews are with zero votes. Many people are agreeing with score 5 reviews.
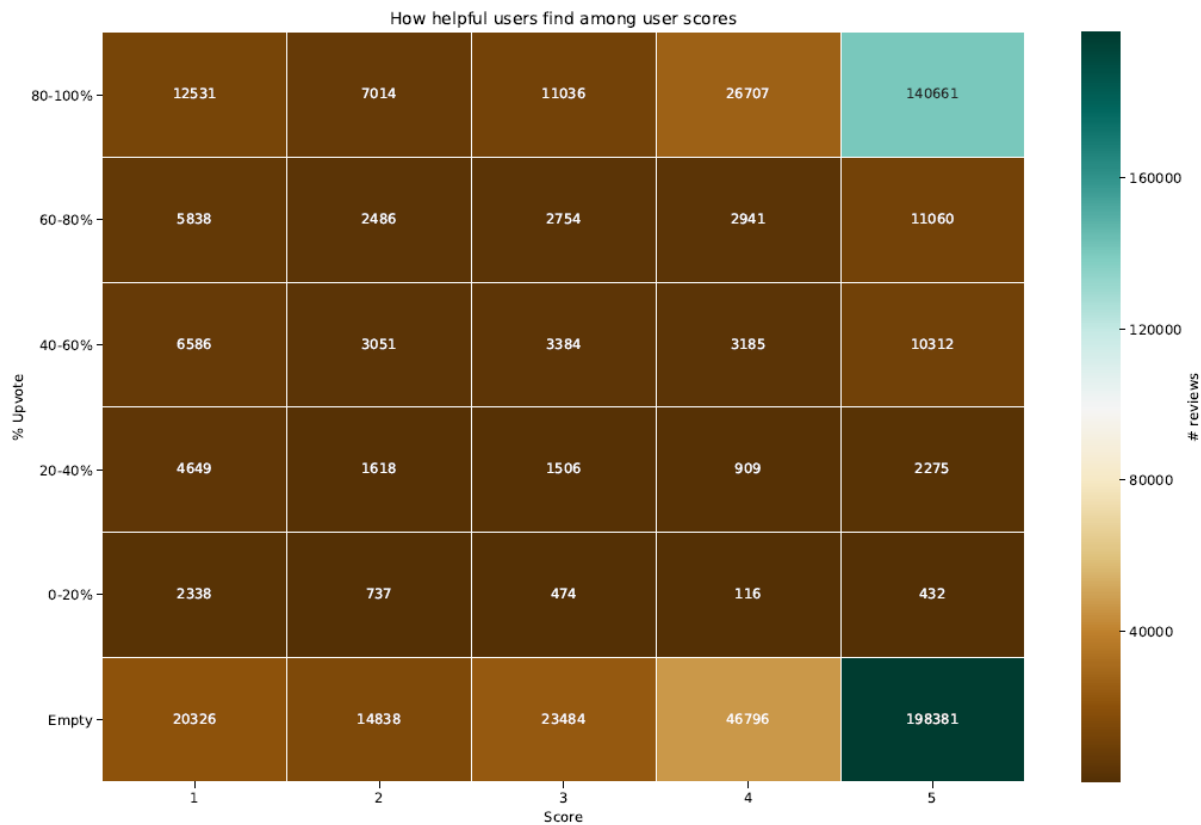
How helpful users find among user scores

| % Upvote | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 80-100% | 12531 | 7014 | 11036 | 26707 | 140661 |
| 60-80% | 5838 | 2486 | 2754 | 2941 | 11060 |
| 40-60% | 6586 | 3051 | 3384 | 3185 | 10312 |
| 20-40% | 4649 | 1618 | 1506 | 909 | 2275 |
| 0-20% | 2338 | 737 | 474 | 116 | 432 |
| Empty | 20326 | 14838 | 23484 | 46796 | 198381 |

Score

**Fig. 3. (f)**

After that we created word cloud, to see the weightage of the words from reviews. Like if customer has used "Disappointed" word many times then the weightage of that word will increase, and that word will look more prominent.  Based on following Wordcloud its very easy to analyze the negative and positive words.

*(Source: (Tutorial) Generate Word Clouds in Python)*

**Sentiment Analysis on Amazon Fine Food Reviews**

Review Score One



Review Score Two

Also while preparing the data, we have used Text mining to remove special characters from the comments, to remove HTML tags, which allowed us to capture words accurately so that we can use those words to build model. Model will not be accurate if it contains special characters as well.

**Sentiment Analysis on Amazon Fine Food Reviews**

# Modeling:

We have used 3 different models, such as Logistic Regression, Decision Tree and Random forest. Also we did run Naïve Bayesian regression but due to memory issue we could not implement it successfully.

**Logistic Regression**: It's a Machine Learning classification algorithm that is used to predict the probability of a categorial dependent variable. In logistic regression, the dependent variable is a binary variable that contains data coded as 1(yes, success, etc.) or 0(no, failure, etc.) In our case, we have categorized customers reviews into 1 & 0 based on the score. Like, score value is more than 3 then it's considered as good review and assigned value as 1 and scores below 3 considered as negative comments, so assigning value to 0. We performed logistic regression on feature Positivity with respect to text columns. We got AUC around 93%, which is good.

*(Source: Logistic Regression In Python)*

### Logistic Regression

```
In [111]:  #Multiple Logistic Regression
           lr = LogisticRegression()
           lr.fit(X_train_vectorized, y_train)
           predictions = lr.predict(vect.transform(X_test))
```
```
C:\Users\Heda\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
C:\Users\Heda\Anaconda3\lib\site-packages\sklearn\svm\base.py:929: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
  "the number of iterations.", ConvergenceWarning)
```
```
In [112]:  print('AUC: ', roc_auc_score(y_test, predictions))

           AUC:  0.9302679895369559
```
```
In [113]:  conf_matr = confusion_matrix(y_test, predictions)
```

**Sentiment Analysis on Amazon Fine Food Reviews**

**Decision Tree**: Decision trees are supervised learning algorithms used for both, classification and regression tasks where we will concentrate on classification in this first part of our decision tree tutorial. We have used "Reviews" on the root node tree and performed Decision tree. We got AUC around 83.42.

*(Source: Decision Tree In Python)*

## Decision Tree

```
In [121]:  #Decision Tree
           dt = tree.DecisionTreeClassifier()
           dt.fit(X_train_vectorized, y_train)
           predictions1 = dt.predict(vect.transform(X_test))

In [122]:  print('AUC: ', roc_auc_score(y_test, predictions1))

           AUC:  0.8342882749002035

In [123]:  conf_matr_dt = confusion_matrix(y_test, predictions1)

In [124]:  sns.heatmap(conf_matr_dt,annot=True,cbar=False, fmt='g')
           plt.ylabel('True Label')
           plt.xlabel('Predicted Label')
           plt.title('Confusion Matrix')

Out[124]:  Text(0.5, 1, 'Confusion Matrix')
```

**Random Forest:** Random forest is a type of supervised machine learning algorithm based on ensemble learning. Ensemble learning is a type of learning where you join different types of algorithms or same algorithm multiple times to form a more powerful prediction model. The random forest algorithm combines multiple algorithm of the same type i.e. multiple decision trees, resulting in a forest of trees, the random forest algorithm can be used for both regression and classification. When we performed Random Forest, we got AUC around 78.69.

*(Source: An Implementation and Explanation of the Random Forest in Python)*

**Sentiment Analysis on Amazon Fine Food Reviews**

## Random Forest ¶

```
In [127]:  #Random Forest
           rf = RandomForestClassifier(n_jobs=-1)
           rf.fit(X_train_vectorized, y_train)
           predictions2 = rf.predict(vect.transform(X_test))

           C:\Users\Heda\Anaconda3\lib\site-packages\sklearn\ensemble\forest.py:245: Futi
           change from 10 in version 0.20 to 100 in 0.22.
             "10 in version 0.20 to 100 in 0.22.", FutureWarning)
```

```
In [128]:  print('AUC: ', roc_auc_score(y_test, predictions2))

           AUC:  0.7869293542832291
```

```
In [129]:  conf_matr_rf = confusion_matrix(y_test, predictions2)
```

Among above three regression models, Logistic regression is better since it gives more accuracy. Apart from this, it does not require too many computational resources, its highly interpretable, it doesn't required input feature to be scaled and its output well-calibrated predicted probabilities. On the other hand Random forest required much more computational resources, owing to the large number of decision trees joined together. Due to their complexity, they require much more time to train than other algorithms.

Also in case of decision tree algorithm, slight change in the data can cause a large change in the structure of the decision tree causes instability. Decision tree often requires more time to train model.

Comparing all these factors, we decided to use Logistic regression.

**Sentiment Analysis on Amazon Fine Food Reviews**

## Evaluation:

### Summary text:

To analyse what kind of word pattern is used by customer on basis of score, we developed word cloud on the summary of review provided by customer. This helped us to understand what are most significant words that is either making the product positive or negative.

### Word Cloud for Score = 1:



Review Score One

From above word cloud we find that words like 'disappointed, worst, waste money, horrible, terrible' denotes that customer are not satisfied with the product and thus rated these products as one. But few words like 'good, great, buy' creates discrepancy when compared to score that is given to product.  So, it is necessary to analyse the text review also in order to make sure if the customer is satisfied or not with the product. It is in future scope of our analysis to compare the score & calculated positivity and make informed decision whether the review tends towards positive or negative side. Because exceptions are to be considered where customer gives score

**Sentiment Analysis on Amazon Fine Food Reviews**

one for something that went wrong like late delivery but like the product. So this approach of comparing score with positivity would help us further enhance the operations of Amazon Fine Foods.

Word Cloud for Score = 2:



Review Score Two

For score two the words seems to tend very slightly towards positive due to words like 'tasty, love, impressed, good' but as score is two, words like 'expensive, old, overpriced' signifies that it is more on negative side.

**Sentiment Analysis on Amazon Fine Food Reviews**

Word Cloud for Score = 3:

Review Score Three

For score three, word cloud resembles that review can be considered as neutral according to scope of our project. There are words equally towards both positive and negative sentiment such as 'good, better, bad, expensive, love, nothing special, great product'. Due to limitation of our project scope we have not considered review score three in text analysis. This reviews can be considered in future to further predict the sentiments accurately.

**Sentiment Analysis on Amazon Fine Food Reviews**

Word Cloud for Score = 4:



Review Score Four

From this word cloud it is evident that score four tends more towards positive sentiment. A conclusion also can be made that customers like gluten free, organic products. Also the customers are satisfied with coffee provided by Amazon Fine Foods.

Word Cloud for Score = 5:



Review Score Five

**Sentiment Analysis on Amazon Fine Food Reviews**

This word cloud shows that customers are satisfied and happy with the products that are been delivered to them. It is also evident that dog food is doing good. Also the customers are liking the taste of the food delivered.

## Machine Learning Models:

We analysed from the word cloud that how the pattern of words tends towards positive sentiment as score goes on increasing. But there are some discrepancies that need to be resolved as discussed for score one wherein we have some positive sentiment words. For this we did text mining and its classification into positivity where '1' stands for positive sentiment and '0' stands for negative sentiment. For text analysis we performed three models viz. logistic regression, decision tree, random forest.

### Logistic Regression:

**AUC** = 0.93 i.e. 93%

**Confusion Matrix:**

**Sentiment Analysis on Amazon Fine Food Reviews**

**Classification Report:**

```
              precision    recall  f1-score   support

           0       0.91      0.88      0.89     20339
           1       0.98      0.98      0.98    111109

    accuracy                           0.97    131448
   macro avg       0.95      0.93      0.94    131448
weighted avg       0.97      0.97      0.97    131448
```

Decision Tree:

**AUC** = 0.834 i.e. 83.4 %

**Confusion Matrix:**



Confusion Matrix

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 14519 | 5820 |
| True 1 | 4954 | 106155 |

**Classification Report:**

```
              precision    recall  f1-score   support

           0       0.75      0.71      0.73     20339
           1       0.95      0.96      0.95    111109

    accuracy                           0.92    131448
   macro avg       0.85      0.83      0.84    131448
weighted avg       0.92      0.92      0.92    131448
```

**Sentiment Analysis on Amazon Fine Food Reviews**

**AUC** = 0.779 i.e. 77.9%

**Confusion Matrix:**



**Classification Report:**

```
              precision    recall  f1-score   support

           0       0.94      0.57      0.71     20339
           1       0.93      0.99      0.96    111109

    accuracy                           0.93    131448
   macro avg       0.93      0.78      0.83    131448
weighted avg       0.93      0.93      0.92    131448
```
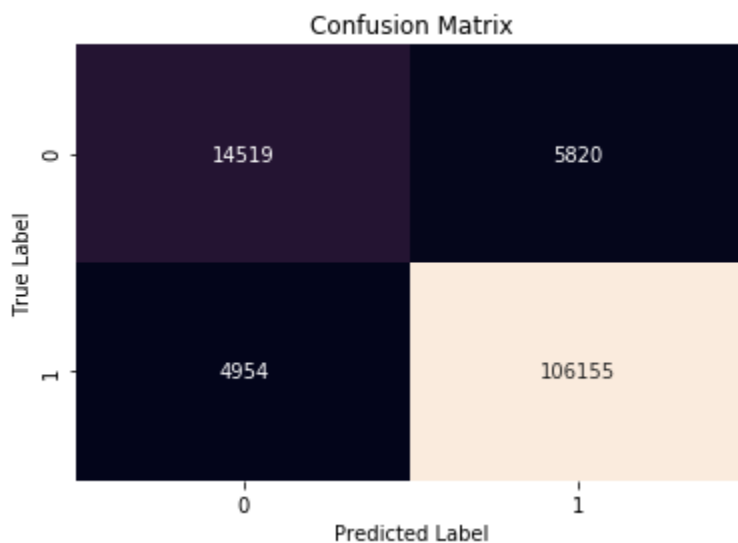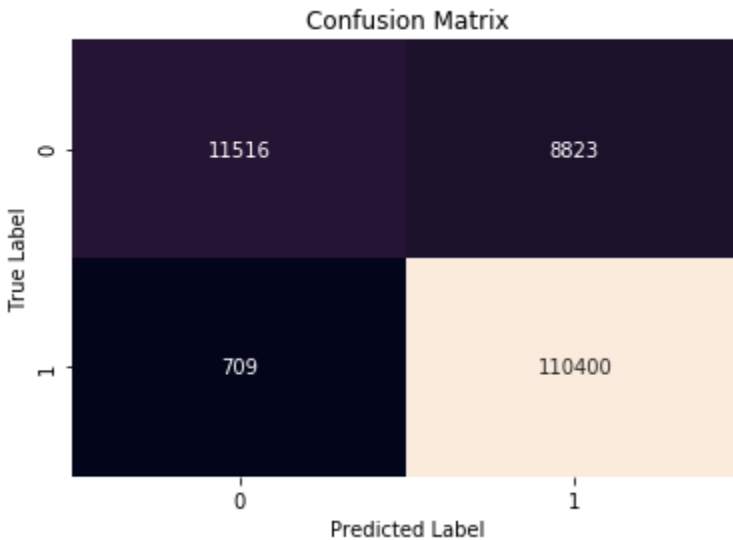
From AUC and classification report we can conclude that logistic regression performs better on classifying the text reviews into positivity i.e. 1 and 0.

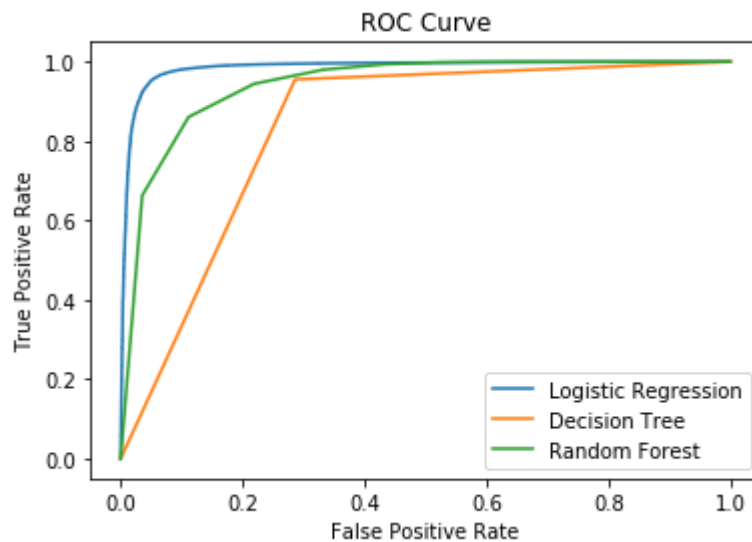Logistic regression and decision tree have better precision and recall for predicting 1's. This may be due to imbalanced data i.e. more reviews are on positive side. While recall for all three models is very low i.e. sensitivity in predicting 0's is less.

Precision for both negative and positive sentiment together with random forest is better as compared to decision tree.

**Sentiment Analysis on Amazon Fine Food Reviews**

F1-score is better for logistic regression in comparison with other models which shows better in precision & recall as well, as f1 score is harmonic mean between precision and recall.

Accuracy is highest for logistic regression followed by random forest and then decision tree.

ROC curve:



It is also evident from ROC curve logistic regression performance is better.

A high threshold value results in a point at bottom left and low threshold value results in a point at top right. This means as we decrease threshold value we get higher TPR at cost of a higher FPR.

*(Source: Understanding Data Science Classification Metrics in Scikit-Learn in Python)*

**Sentiment Analysis on Amazon Fine Food Reviews**

# Deployment:

So far, we have cleaned the data, prepared it, performed regression models by dividing the data into test and train to analyze whether the reviews are positive or negative, performed regression model on word count, to identify which words have occurred the most. We decided to go with Logistic regression model, since it provided great accuracy.

After performing all these, the analysis carried out to focus on one specific user, on what he/she likes in term of fine food, based on the reviews he/she had given in the past. This can be expanded to all the users later on. We performed this, so that Amazon can focus on the customers, they can analyze the customers behavior like what he/she likes, accordingly Amazon can recommend the relevant food. Amazon can do customer-oriented marketing. As a result it will help Amazon to grow their customers. Also this approach will save food and money as Amazon will deliver only those products which customers likes the most. It will also generate revenue with minimum investment. For more clarity, lets analyze the following graphs. In the Fig. 6. (a) , we have identified the users who has given most positive reviews based on mean of hi all the score. We have also identified the number of reviews he has made.

```
                                            Score count   Score mean
UserId          ProfileName
A3OXHLG6DIBRW8  C. F. Hill "CFH"                   448      4.535714
A1YUL9PCJR3JTY  O. Brown "Ms. O. Khannah-Brown"    421      4.494062
AY12DBB0U420B   Gary Peterson                      389      4.647815
A281NPSIMI1C2R  Rebecca of Amazon "The Rebecca Review"  365  4.841096
A1Z54EM24Y40LL  c2                                 256      4.453125
A1TMAVN4CEM8U8  Gunner                             204      4.833333
A2MUGFV2TDQ47K  Lynrie "Oh HELL no"                201      3.751244
A3TVZM3ZIXG8YW  christopher hayes                  199      1.000000
A3PJZ8TU8FDQ1K  Jared Castle                       178      4.601124
AQQLWCMRNDFGI   Steven A. Peterson                 176      3.954545
```

**Sentiment Analysis on Amazon Fine Food Reviews**

In the Fig. 6. (b), we have focused on only user that is C.F. Hill and we identified his number of ratings. From below graph we can he has given 5 ratings to almost 250 products and less than 10 he has given 1 rating. We will do the same for all the users.

**Fig.6. (a)**



Score distribution of user C. F. Hill "CFH" review

**Fig. 6. (b)**
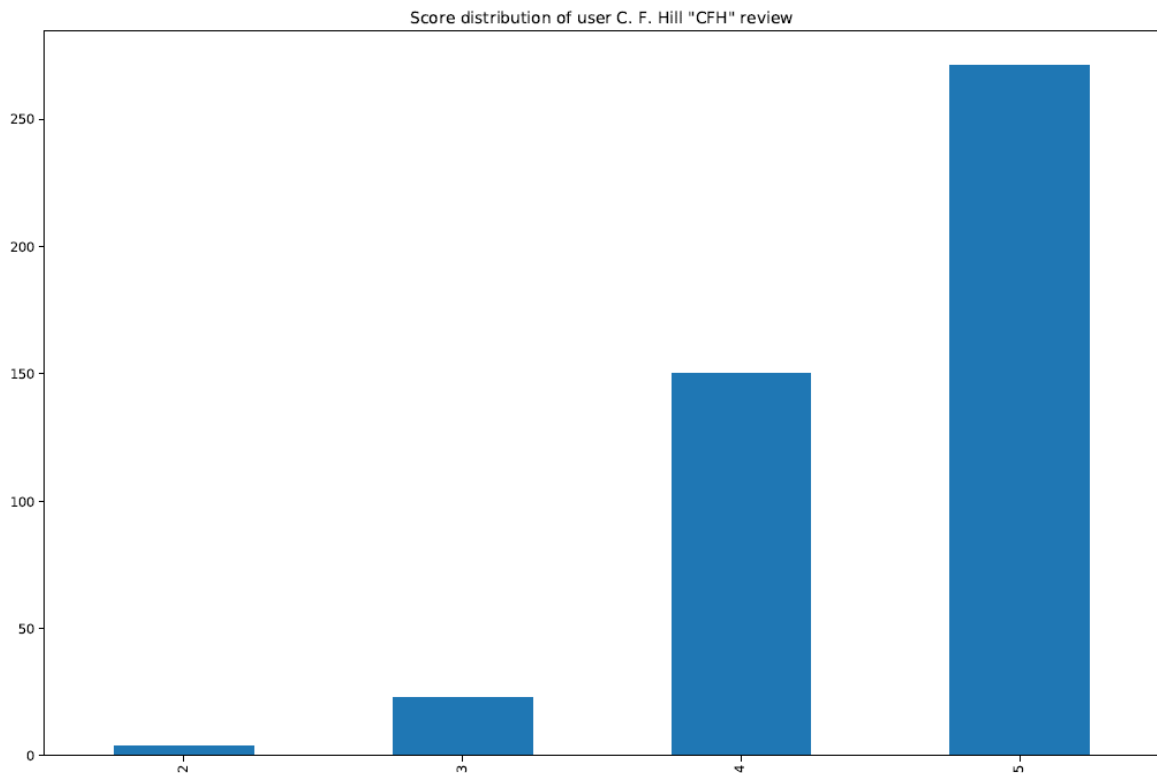
So, based on the above graph it is helpful to Amazon to focus more on the products for which C.F. Hill has given 4 or 5 ratings. So that they will deliver such products more, as a result it increase their revenue.

Amazon could also ask for feedback from the customer who has given below 3 rating so that they can improve their products

**Sentiment Analysis on Amazon Fine Food Reviews**

As we will be deploying Logistic regression, the issues we need to take into consideration are it will not perform well with independent variables that are not correlated to the target variable and are very similar or correlated to each other. So we will have to make sure that we will remove variable that are less correlated with target variable. Sometimes its vulnerable to overfitting.

The only risk currently we found that the model is not 100% accurate. So suppose if we use dataset which is less than 100 then the model might not identify the positive or negative comment accurately. Sometimes it will give wrong confusion matrix. So we working on building a model or methods which will also work for small dataset and model will be more accurate.

**Sentiment Analysis on Amazon Fine Food Reviews**

# References:

1. SNAP: Web data: Amazon Fine Foods reviews

   **SNAP: Web data: Amazon Fine Foods reviews.** *(2019). Snap.stanford.edu.*

2. Logistic Regression In Python

   **Logistic Regression In Python.** *(2019). Medium.*
   *https://towardsdatascience.com/logistic-regression-python-7c451928efee*

3. Decision Tree In Python

   **Decision Tree In Python.** *(2019). Medium.*
   *https://towardsdatascience.com/decision-tree-in-python-b433ae57fb93*

4. Decision Trees in Python with Scikit-Learn

   **Decision Trees in Python with Scikit-Learn.** *(2018). Stack Abuse.*
   *https://stackabuse.com/decision-trees-in-python-with-scikit-learn/*

5. An Implementation and Explanation of the Random Forest in Python

   **An Implementation and Explanation of the Random Forest in Python.** *(2018). Medium.*
   *https://towardsdatascience.com/an-implementation-and-explanation-of-the-random-forest-in-python-77bf308a9b76*

6. 3.2.4.3.1. sklearn.ensemble.RandomForestClassifier — scikit-learn 0.21.3 documentation

   **3.2.4.3.1. sklearn.ensemble.RandomForestClassifier — scikit-learn 0.21.3 documentation.** *(2019). Scikit-learn.org.*
   *from https://scikit-learn.org/stable/modules/generated/sklearn*

7. Understanding Data Science Classification Metrics in Scikit-Learn in Python

   **Understanding Data Science Classification Metrics in Scikit-Learn in Python.** *(2019). Medium.*
   *https://towardsdatascience.com/understanding-data-science-classification-metrics-in-scikit-learn-in-python-3bc336865019*

**Sentiment Analysis on Amazon Fine Food Reviews**

8.  (Tutorial) Generate Word Clouds in Python

    **(Tutorial) Generate Word Clouds in Python.** *(2019). DataCamp Community.*
    *https://www.datacamp.com/community/tutorials/wordcloud-python*


9.  Text Analytics for Beginners using NLTK

    **Text Analytics for Beginners using NLTK.** *(2019). DataCamp Community.*
    *https://www.datacamp.com/community/tutorials/text-analytics-beginners-nltk*

**Sentiment Analysis on Amazon Fine Food Reviews**

# Project timeline and member contribution:



| | | Task Mode | Task Name | Duration | Start | Finish | Predecessors |
|---|---|---|---|---|---|---|---|
| 1 | ✓ | 📌 | Brainstorming | 5 days | Mon 16-09-19 | Fri 20-09-19 | |
| 2 | ✓ | 📌 | Data Gathering | 2 days | Mon 23-09-19 | Tue 24-09-19 | 1 |
| 3 | ✓ | 📌 | Project Proposal | 3 days | Wed 25-09-19 | Fri 27-09-19 | 1,2 |
| 4 | ✓ | 📌 | Data Exploration | 5 days | Mon 30-09-19 | Fri 04-10-19 | 2,3 |
| 5 | ✓ | 📌 | Data Cleaning | 5 days | Mon 07-10-19 | Fri 11-10-19 | 2,4 |
| 6 | ✓ | 📌 | Data Preparation | 7 days | Mon 14-10-19 | Tue 22-10-19 | 2,4,5 |
| 7 | ✓ | 📌 | Logistic Regression (Score Prediction) & other models | 5 days | Wed 23-10-19 | Tue 29-10-19 | 2,4,5,6 |
| 8 | ✓ | 📌 | Extracting Features from text data | 10 days | Wed 23-10-19 | Tue 05-11-19 | 2,4,5,6 |
| 9 | ✓ | 📌 | Result Evaluation | 4 days | Wed 06-11-19 | Mon 11-11-19 | 7,8 |
| 10 | ✓ | 📌 | Tuning the results if needed | 5 days | Tue 12-11-19 | Mon 18-11-19 | 7,8,9 |
| 11 | | 📌 | Documentation | 8 days | Tue 19-11-19 | Thu 28-11-19 | 2,4,5,6,7,8,9,10 |
| 12 | | 📌 | Project close out & Presentaion | 1 day | Mon 02-12-19 | Mon 02-12-19 | |

**Sentiment Analysis on Amazon Fine Food Reviews**