# Applied ML Coursework

Sohel Arsad Salam
*COMP-1804-M01-2021-22 Applied Machine Learning*
001166291

*Abstract*—**This document provides a basic idea of executing ML solutions for a particular issue utilizing the data-set relegated. Planning an AI arrangement requires thinking about a few parts of the issue, the accessibility and neatness of information and relating explanations, nature of the issue tended to, technique decision, assessment among others.As online commercial centers have been famous during the previous many years, the online dealers and vendors request that their buyers share their viewpoints about the items they have purchased.Audits on Amazon are not just connected with the item yet in addition the assistance given to the clients. If clients get clear bifurcation about item audits and administration audits it will be simpler for them to take the choice, in this paper we propose a framework that plays out the arrangement of client surveys followed by tracking down feeling of the audits. A standard based extraction of item include feeling is likewise finished. Likewise we give a perception to our outcome rundown. Examining this huge measure of sentiments is likewise hard and tedious for item makers. This proposition considers the issue of arranging surveys by their generally semantic (positive or negative). To lead the review, two different supervised ML procedures, Logistic Regression Model and Random Forest, has been endeavored on items from Amazon. Their exactness have then been analyzed. The outcomes showed that the Logistic Regression Model approach beats the Random Forest approach when the informational index is greater. In any case, the two calculations arrived at promising exactness of at minimum 90 percent.**

## I. Introduction and related work

Amazon is one of the biggest web-based seller in the World. Individuals frequently look over the items and surveys of the item prior to purchasing the item on amazon itself. This report discusses an arrangement issue which is taking in client audits of various items (i.e a computer game furthermore, an instrument) and the evaluations given to them. This report discusses an order issue which is taking in client audits of various items (i.e a computer game furthermore, an instrument) and the appraisals given to them. The dataset contains two objective factors one being the kind of item and the other one being the appraisals of the clients which goes from - 1 to 5. This model will help the amazon laborers to remarks without evaluations and fill their appraisals and if the rating is poor , it very well may be utilized for investigation and work on the said result of the seller. This is a difficult issue on the grounds that since the surveys are text information , when they are tokenised and switched over completely to highlights , the elements are more than 80,000 and the most serious issue is that the information is intensely imbalanced. Out of the 6 classes in the audit score of the item , the greater part

of the information were concentrated around the class 4,5 which intends that the items has for the most part just got positive audits yet it would make the classifier one-sided to the positive remarks and it may be anticipating the positive classes. since the issue within reach is generally delicate for negative items since that is the use case we are attempting to tackle for. The dataset contains audit id ,audit score , confirmed status , survey and item classification. where survey score and the item class are the objective factors. The preprocessing for text information, for example, eliminating accentuation , lemmatisation , stemming ,one hot encoding to convert text classes to mathematical qualities, n-grams to take single twofold mix of words and eliminating stopwords were utilized prior to vectorising the information utilizing TFIDF vectoriser also, SMOTE was utilized to adjust the information since the multiclass target was intensely imbalanced. The calculations used to order the information are Logistic Regression and Random Forest since this is an order issue these are the classifiers utilized. Logistic Regression utilizes the sigmoid capacity and returns probability of the expectations which are switched over completely to classes , while the Random Forest utilizes numerous choice trees to initialise various trees on irregular haphazardly blend of highlights and returns the typical likelihood to group a datapoint.

## II. Ethical discussion

In under an age, the unstable development of AI (ML) has turned into an extraordinary social, political, and monetary power the world over. No matter what, the inescapability of arranged figuring, of computerized interconnectedness, and of universal information extraction along with fast progress made in figuring power and algorithmic methods are currently introducing partners across genuinely every area of society, and at each financial level, with extraordinary open doors as well as huge difficulties. The open doors might well appear to be unlimited. Datadriven bits of knowledge created by ML advancements have currently began to progress vital elements of human prosperity and further developed possibilities for a more reasonable future. In the field of medical services, for example, biomedical ML applications are permitting specialists to more readily target malignant growth drugs, to distinguish sicknesses prior and all the more really, and to do surgeries with uncommon accuracy. In numerous other future-basic fields as well, from natural science to energy the board, ML applications are combatting environmental

change and deforestation, supporting biodiversity, catalyzing rural efficiency, creating 'more brilliant,' more effective urban areas, and assisting with giving additional opportunities for the democratized appropriation of fundamental merchandise also, administrations. As these innovatively initiated gambles have come into more clear view, pundits have started to voice authentic worries about the threats to individual furthermore, social opportunity presented by the quick multiplication of always computationally strong applications of ML. Some have zeroed in on the 'huge tech' driven political economy of observation free enterprise wherein the uncovered double-dealing of standards of conduct serves as an apparatus for buyer control and corporate benefit. Others have seen how boundless government utilization of advanced following and computerization upheld independent direction can — when completed carelessly, flippantly, and without comprehensive local area inclusion — capacity to build up firmly established examples of destitution, disparity, also, underestimation.

## III. DATASET PREPARATION

The initial step for leading the examination incorporates information assortment for preparing also, testing the classifiers. The document was changed over completely to the Comma Isolated Values (CSV) design, as it is more helpful for python to deal with this sort of records. The informational collection comprises of multiple surveys of different excellence products. For setting up the ideal information a basic code was written in python to eliminate the pointless highlights.

Eliminate Accentuations : This is utilized to eliminate the punc- tuations present in the information involving the accentuations in string library. They eliminate every one of the accentuations in the string.

Eliminate Stopwords : This is utilized to eliminate the undesirable rehashing words in the english language that doesn't add any importance to the audits. the stopwords bundle was introduced from nltk. they were utilized as a kind of perspective to eliminate from the information.

Tokenisation : This is the most common way of separating text information , the audits into endlessly sentences into word level granularity.

Unigrams and Bigrams : From the tokens accessible in each audit a blend of unigrams (1 word mixes) and bi-grams (2 progressive word blends) are recovered and they structure the highlights for the model.

Lemmatisation : This is the most common way of changing over each words into the base lemma. they convert the words into their root words no matter what their strained.

Stemming : This is the method involved with changing over each word into its base structure without the endings , this is more inclined to distinction since the strained of the word causes changes. the equivalent word in its past tense and future tense might have an alternate type of stemming word.

Class Balancing : Class adjusting is finished utilizing calculations like Destroyed for upsampling and downsampling and ran- dom inspecting to adjust the multiclass target factors. The classes were intensely imbalanced so they were upsampled utilizing Destroyed which gave the best precision contrasted with the remainder of the calculations.

Vectorisation : Vectorisation is the method involved with changing over text into mathematical elements since the AI models cant peruse text information , they must be changed over. The vectorisation technique utilized here is TF-IDF. Term Recurrence and Opposite Archive Recurrence which loads words in view of the less occurence in the survey , contrasted with the words that are most regularly occuring.

Splitting of Data : Subsequent to doing all the previously mentioned preprocessing the information is parted into three sets , the train , test and the approval datasets. The train is utilized to prepare the model and the test is used to evaluate its exactness and the approval is for the model to anticipate and convey realtime results or concealed information later assessing the model.

Many highlights were eliminated with the exception of the rundown of the audit, the text of the actual survey, score and productId. In the main examination the entire informational index was utilized. Since the quantity of surveys were sufficient to obtain a sensible outcome from the classifiers the audits with three stars were discarded to keep away from any entanglement while preparing the algorithms. However, in the second analysis because of the modest number of information the audits with three stars were additionally thought to be as negative. A similar code was then used to name the information.

## IV. METHODS

There are two models required here in the forecast of two objective factors.

- 1st Methodology (Logistic Regression Method):
  - For target variable 1 (Multi class arrangement for audit score going from - 1 to 5) . In the wake of doing the preprocessing the data.(removing accentuations , stopwords, tokenisation , stemming, lemmatisation and vectorisation) the information is taken care of into Logistic Regression model. The information was checked for irregular characteristics , the information was vigorously imbalanced. the classes were - 1,1,2,3,4,5. The class 5 had 18704 , class 4 had 6015 , class 3 had 3038 , class 2 had 1512 and class 1 had 1820 and the class - 1 has 904 lines. after applying SMOTE to upsample all classes were adjusted. the calculated model returned an exactness of 51 percent. Since there were various classes the models that work better were Random Forest since they make numerous trees and foresee , Logistic Regression was shown to comparatively be the best indicator.
- 2nd Methodology (Random Forest Method):
  - For target variable 2 (Multi class arrangement for audit score going from - 1 to 5) . In the wake of doing the preprocessing the data.(removing accentuations , stopwords, tokenisation , stemming, lemmatisation and vectorisation) the information is taken care of

into Random Forest model. The data was checked for sporadic attributes , the data was overwhelmingly imbalanced. the classes were - 1,1,2,3,4,5. The class 5 had 18804 , class 4 had 6215 , class 3 had 3138 , class 2 had 1472 and class 1 had 1880 and the class -1 has 984 lines. in the wake of applying Destroyed to upsample all classes were changed. the determined model returned a precision of 61 percent.

## V. EXPERIMENTS

- Model - 1 : Logistic Regression Algorithm :
  - Logistic Regression Algorithm is one the most ordinarily utilized arrangement calculations , It is a direct relapse applied on information on which the sigmoid capacity is applied to change the expectations over completely to its probablities , and the straight line of the straight relapse is changed over into a S molded bended line. the information is cleaned , eliminated of accentuations , stopwords , tokenised , ngrams had been applied so mix of single , twofold and triple progression mix of words are additionally recovered and afterward they are vectorised utilizing the tf-idf vectoriser the paramters were utilized, for example, mindf and a different scope of values were given from 2 till 10 , there wasnt any significant changes in the vectors recovered. then the information was checked for lopsided characteristics in view of the objective variable-2 which was the survey score , the audit score had 6 classes going from - 1 ,1,2,3,4,5 the least signifying the most unfortunate and 5 being the best evaluations accommodated the item. Expanding PCA alongside customary PCA and shortened SVD were additionally carried out to check whether the aspect decrease assists with the order of the multi classes , however there wasnt any improvement in the class isolation given the parts in the scope of 2,5,10. The information was viewed as vigorously imbalanced the audit score 5 being the most elevated of around 18,000 surveys and - 1 audit score had barely 150 surveys and each other had around 1000 surveys. this was adjusted utilizing destroyed upon the vectorised information. This information was parted into train and test. The train information alongside its objective qualities were fitted on the model and afterward tried on the test information by trasnforming the test information into the comparative vectors and when anticipated the model had 47 rate exactness. The PCA didnt assist with supporting the exactness.
- Model - 2 : Random Forest Algorithm :
  - Random Forest is a gathering of order based trees that haven't been managed. It performs commendably in an assortment of certifiable issues since it is unaffected by commotion in the dataset and the peril of overfitting is negligible. It is faster than numerous

other tree-based calculations and upgrades precision for testing and approval information. The total of the forecasts of individual choice tree calculations is known as Random Forest. While making an irregular tree, there are a few choices for tuning the Random Forest's presentation. The strength of the Random Forest classifier comes from its arrangement trees. Since it is framed from "irregular" subsets of the information, and the eventual outcome is contrasted with different trees that have additionally been shaped "haphazardly", the calculation watches well against "overfitting" from boisterous information focuses that might have more effect on a solitary choice tree calculation. The "irregular" development of the trees guarantees that there is minimal opportunity for a solid predisposition to be available in the information during tree development. Random Forest additionally function admirably with high layered information, as on account of our TFIDF highlights. Since Random Forest functions admirably with high layered information, as is serious with different calculations, for example, SVM without the high preparation cost it is chosen for our review.

### A. Evaluation

Evaluation metrics for classification algorithms:

- Precision :
  - Precision is the quantity of appropriate grouping of up-sides which are distinguished by the model on the whole dataset's sub-populace of up-sides. It is fundamentally the way that well the model can distinguish the up-sides.

$$PRECISION = \frac{true_Positive}{true_Positive + false_Positive} \quad (1)$$

- Recall
  - Recall is the quantity of genuine up-sides appropriately distinguished by the model , which is the quantity of genuine up-sides for complete arrangement of genuine up-sides and misleading negatives.

$$RECALL = \frac{true_Positive}{true_Positive + false_Negative} \quad (2)$$

- F1-score
  - F1- score tells how good the model is performing both as far as accuracy and review as opposed to picking one measurement over the other

$$F1score = 2 * \frac{PRECISION * RECALL}{PRECISION + RECALL} \quad (3)$$

Evaluation metrics values for Logistic Regression were :
Evaluation metrics values for Method -2 were :

TABLE I
E

valuation metrics for logistic model for target variable -1 labels prediction

| LogisticRegression | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.98 | 0.93 | 0.95 | 4352 |
| 1 | 0.86 | 0.95 | 0.91 | 2048 |

TABLE II
E

valuation metrics for Random Forest model for target variable -2 labels

prediction

| RandomForest | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| 0 | 0.97 | 0.93 | 0.95 | 4352 |
| 1 | 0.87 | 0.95 | 0.91 | 2048 |

TABLE III
T

he table shows the evaluation metrics for naive bayes model for target variable -2 labels prediction

| LogisticRegression | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| -1 | 0.00 | 0.00 | 0.00 | 164 |
| 1 | 0.61 | 0.49 | 0.54 | 381 |
| 2 | 0.45 | 0.18 | 0.26 | 351 |
| 3 | 0.43 | 0.27 | 0.33 | 601 |
| 4 | 0.46 | 0.33 | 0.39 | 1262 |
| 5 | 0.73 | 0.91 | 0.81 | 3822 |

items. The grouping of weather the item was an instrument or computer game had received the most noteworthy arrangement exactness and the multi class order for the survey score didn't earn a lot of precision as there were such a large number of classes to be anticipated and strategies like PCA and Expanding PCA additionally didn't perform well as the quantity of parts went high the machine couldn't handle the heap , assuming that the classes were binned together into 3 or 2 sub classes , the precision might have been something more.

## VI. DISCUSSION AND FUTURE WORK

The primary objective of this study was to figure out which AI calculation of Random Forest or Logistic Regression techniques performs better in the assignment of text characterization. This was achieved by involving the Amazon reviews dataset as informational collection. The classifiers were assessed by looking at their exactnesses in different cases of trials. The outcomes from the main arrangement of tests demonstrated that the Logistic Regression approach got preferable precision over the Random Forest in the two situations where the calculations applied on the audits and whenever they have been applied on the outlines. The difference in exactnesses between these methodologies is anyway tiny. From this trial it tends to be figured out that thoroughly prepared AI calculations with enough information as preparing informational index can perform awesome grouping. As far as correctnesses, Logistic Regression will in general show improvement over Random Forest, albeit the differences aren't exceptionally enormous, and the calculations can arrive at more than 90 percent of order accurately. Another explanation can be the audits with 3 stars which for the most part are arranged as impartial. Notwithstanding, In the second analysis because of the little size of information set they were thought of as negative. This could have affected the consequence of the second trial. In future review ordering them would intrigue as sure to see whether it gives a different result.

## VII. CONCLUSIONS

This report closes with the outcomes that were acquired with the analyses performed on the amazon surveys on the two