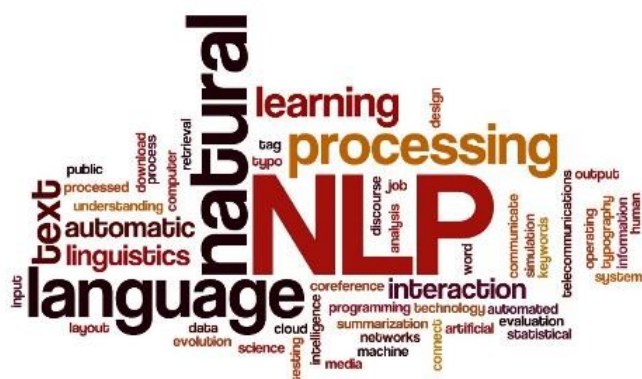


Natural Language Processing (NLP) Application Capstone Project

April 11, 2020



Automated Ticket Assignment



Program	Post Graduate Program (PGP) in Artificial Intelligence and Machine Learning(AIML)
Batch	April 2019 – April 2020
Group	3
Project Guide	Sanjay Tiwary
Program Manager	Anurag Shah
Team Members	Gaurav Walia, Karishma Dcosta, Lavanya Harry Pandian, Pallavi Kumari, Swati Tyagi
Deliverable	Final Report
Submission Date	19th April 2020

1 Project Goal

One of the key activities of any IT function is to ensure there is no impact to the Business operations through Incident Management process. An incident is an unplanned interruption to an IT service or reduction in the quality of an IT service that affects the Users and the Business. The main goal of the Incident Management process is to provide a quick fix / workarounds or solutions that resolves the interruption and restores the service to its full capacity to ensure no business impact. These incidents are recorded as tickets that are created by various stakeholders (Business Users, IT Users and Monitoring Tools) within IT Service Management Tool and are assigned to Service Desk teams (L1 / L2 teams). The goal of this project is to build a classifier that can classify the tickets by analyzing the text using Natural Language Processing(NLP) techniques in AIML.

2 Summary of problem statement, data and findings

In this section, we describe the problem statement: explaining the current situation, opportunities for improvement and data findings: data requirement, size and source of data with its challenges and techniques to overcome the same.

2.1 Problem Statement

Current Situation

Given the data that is collected from the IT Service Management Tool, the issues are recorded as tickets and are assigned to respective groups based on the type of issues that need to be addressed. Assigning the incidents to the appropriate group has critical importance to provide improved user satisfaction while ensuring better allocation of support resources, thus maintaining the organization's efficiency in the service. However, the assignment of incidents to appropriate IT groups is still a manual process in many of the IT organizations. Manual assignment of incidents is time consuming and requires human efforts. There may be mistakes due to human errors and resource consumption is carried out ineffectively because of the misaddressing. On the other hand, manual assignment increases the response and resolution times which result in user satisfaction deterioration or poor customer service.

Opportunity for improvement through ML

This manual process can be improvised using machine learning based systems such as automatic ticket classification mechanisms that would:

- Reduce or remove the count for human error
- Use the allocated resources efficiently
- Ensure efficient ticket classification
- Provide quick solutions and turn-around times for the organization

By leveraging the AI technology, we shall build a classifier that can classify the tickets into respective Groups by analyzing the text using NLP techniques in AIML.

2.2 Data Findings

In Natural Language Processing (NLP), most of the data in the form of documents and text contain many words that are redundant for text classification, such as stop-words, misspellings, slangs; and contain various languages since the users could potentially be located globally. **Data Requirement** To understand the tickets, we require the past ticket information comprising of the ticket summary/ title which captures the essence of the issue, the detailed description for additional details, the user information and the group assigned to the respective tickets. Additional information such as separate fields for timestamp, geographic location of user, etc., would be useful in understanding the traffic and geo of the tickets logged to assign resources as per the demand of the tickets. The Dataset used for the project can be referred from the following location in an excel format(*.xlsx): <https://drive.google.com/drive/u/0/folders/1xOCdNI2R5hiodskIJbj-QySMQs6ccehL>

Source of data and challenges

The data that has been captured by the IT Management System Tools is unclean. The data requires to be devoid from noise, punctuations, misspellings, htmls, etc. as a start. Further, data should be pre-processed to remove words which do not contribute to context (stop-words) and extract meaningful words (tokens) to feed the data to modelling algorithms.

	Short_description	Description	Caller	Group
0	login issue	-verified user details.(employee# & manager na...	spxjnrwir pjlcqds	GRP_0
1	outlook	\r\n\r\nreceived from: hmjdrvpb.komuaywn@gmail...	hmjdrvpb komuaywn	GRP_0
2	cant log in to vpn	\r\n\r\nreceived from: eylqgodm.ybqkwiam@gmail...	eylqgodm ybqkwiam	GRP_0
3	unable to access hr_tool page	unable to access hr_tool page	xbkucsvz gcpydteq	GRP_0
4	skype error	skype error	owlgqjme qhcozdfx	GRP_0

Figure 1: Raw data read from the Excel file with 4 columns

The methods employed to clean our data are used from the NLTK and Regular Expression (re) library.

Size of the data

Duplicate entries - The dataset consists of 8500 entries of tickets. On analyzing for duplicate entries across all the 4 columns, 83 duplicates were observed and removed thus leading to unique 8417 values in the dataset.

Class Imbalance

Datasets require proper representation of Class information, i.e equal representation of all Groups. This would enable the modelling algorithms to be trained on equal amounts of data of any given

```
# Drop duplicate rows
df_v1 = df
df_v1 = df_v1.drop_duplicates(keep='first', inplace=False)
df_v1.shape

(8417, 4)
```

Figure 2: Dropping duplicate entries in the dataframe

class (Group). However, this is not the case, and we observe data imbalance in the dataset. Given the Group information, the Unique Group Count = 74. However, there is a class imbalance w.r.t the representation of the groups in the data. Out of the total 8500 tickets, about 47% represents Group_0 tickets. One way to control the group data and maintain imbalance is by setting a 'Threshold' value which filters out the minority Group data. The threshold value is set at default 50.

```
# Reset Assignment Group for group types with less data
Frequency_Threshold = 50
count = df_v1['Group'].value_counts(ascending=True)
idx = count[count.lt(Frequency_Threshold)].index
df_v1.loc[df_v1['Group'].isin(idx), 'Group'] = 'GRP_Manual'
print("Updated unique group types", df_v1['Group'].nunique())
df_v1['Group'].value_counts(ascending=True)
```

Figure 3: Setting a threshold at 50 to control the class imbalance during modelling

We can tune this Threshold value based on the business requirements to filter the appropriate information into our dataset. However, there are drawbacks to this method which urges us to focus on sampling techniques. Sampling techniques would enable us to down sample the majority classes or/and upsample the minority classes.

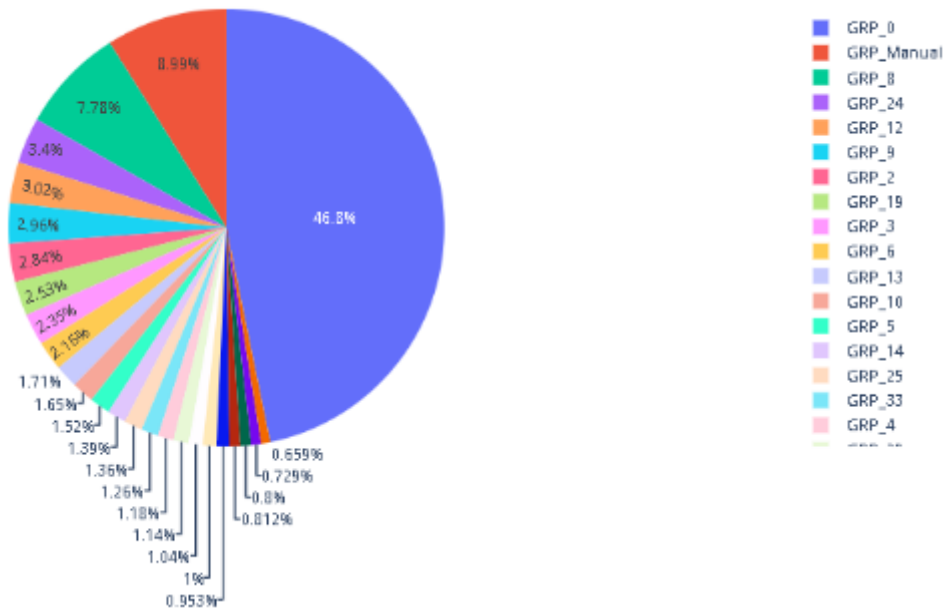


Figure 4: Out of the total 8500 tickets, 47% represents Group_0 tickets

```
# Reset Assignment Group for group types with less data
Frequency_Threshold = 5 #50
count = df_v1['Group'].value_counts(ascending=True)
idx = count[count.lt(Frequency_Threshold)].index
df_v1.loc[df_v1['Group'].isin(idx), 'Group'] = 'GRP_Manual'
print("Updated unique group types",df_v1['Group'].nunique())
df_v1['Group'].value_counts(ascending=True)
```

Figure 5: Setting a threshold=5 to control the class imbalance during modelling

3 Summary of the Approach to EDA and Pre-processing

3.1 Analyze and understand the structure of data

Reading the dataset

The dataset is located in the google drive and accessed using Pandas library. This file is then stored in a dataframe. While reading the file, 'Assignment group' is renamed to 'Group' and 'Short description' to 'Short_description' as given below.

Viewing dataframe with dataframe.head()

We analyze the shape of the data to get a basic understanding of the data and features by using dataframe.shape and dataframe.describe() respectively.

```
# Read Dataset
file_name = "Ticket_Data.xlsx"
df = pd.read_excel(file_name, encoding='cp1252')
df = df.rename(columns = {"Short description": "Short_description",
                          "Assignment group": "Group"})
```

Figure 6: Reading the dataset using Pandas library

df.head()

	Short_description	Description	Caller	Group
0	login issue	-verified user details.(employee# & manager na...	spxjnwir pjlcqds	GRP_0
1	outlook	\n\nreceived from: hmjdrvpb.komuaywn@gmail...	hmjdrvpb komuaywn	GRP_0
2	cant log in to vpn	\n\nreceived from: eylqgodm.ybqkwiam@gmail...	eylqgodm ybqkwiam	GRP_0
3	unable to access hr_tool page	unable to access hr_tool page	xbkucsvz gcpydteq	GRP_0
4	skype error	skype error	owlgqjme qhcozdfx	GRP_0

Figure 7: Getting a preview of the dataframe

- There are 4 columns namely Short_description, Description, Caller and Group. The total number of entries are 8500 in the dataframe.
- The count is different for each column, indicating missing values at first glimpse.
- Unique denotes the unique values in each column Eg. There are 74 unique groups in the dataframe.
- Displays the top word or content in the columns
- Frequency captures the frequency at which the top word/content appears in the columns.

```
# Checking Shape of the data
print("Data shape:", df.shape)
print("Data Description:")
df.describe()
```

Data shape: (8500, 4)

Data Description:

	Short_description	Description	Caller	Group
count	8492	8499	8500	8500
unique	7481	7817	2950	74
top	password reset	the	bpctwhsn kzqsbmtp	GRP_0
freq	38	56	810	3976

Figure 8: Checking shape and basic description of the data

Checking for Missing values

On analyzing the raw data, the fields (Short Description, Description, Caller and Group) were verified for missing values. It was observed that there are 8 missing values in Column - Short Description and 1 missing value in Column - Description and none in Column - Group. The missing values were addressed by imputing a stop word 'the', which will be processed in the stop-word removal or text cleaning.

```
# Check for number of null values in each columns
print("Total Null Values in data:", df_v1.isnull().sum().sum())
print("\nNull Values accross columns:\n", df_v1.isnull().sum())
print("\nData with 'Null' Short Description")
df_v1.loc[df_v1['Short_description'].isnull()==True]
```

Total Null Values in data: 9

Null Values accross columns:

Short_description	8
Description	1
Group	0

Figure 9: Checking for missing values and imputing data in the dataset columns

3.2 Visualize Data

Visuals are a great way to analyze data and get an idea about the data that we are handling. There are many packages in python that help visualize data. We have chosen Word Clouds and defined a function to analyze:

- Most frequent words in raw Short_description
- Most frequent words in raw Description
- Clean data in Summary field

```
[ ] 1 def wordCloudText(data, title):
2     title = ("Most Frequent words in ") + title
3     stopwords = set(STOPWORDS)
4     wordcloud = WordCloud(background_color='black', stopwords=stopwords, max_words=200,
5                           max_font_size=40, random_state=42).generate(str(data))
6     print(wordcloud)
7     fig = plt.figure(1,figsize = (20, 8))
8     plt.imshow(wordcloud)
9     plt.title(title ,fontsize=30)
10    plt.axis('off')
11    plt.show()
```

Figure 10: Word Cloud function

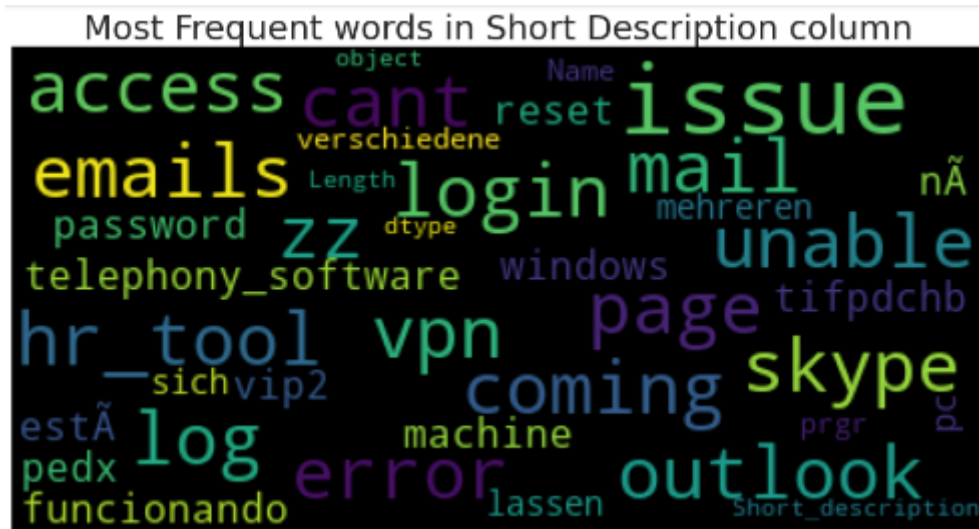


Figure 11: Most frequent words in Short_descriptionbeforecleaningdata

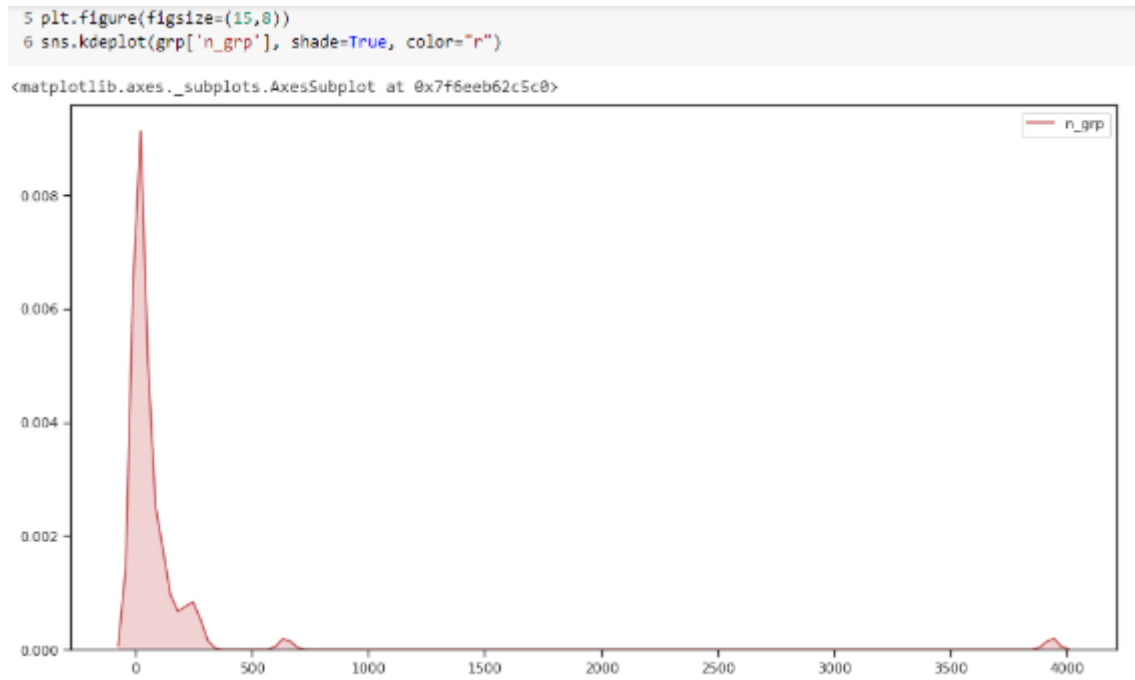


Figure 15: kde representing the frequency distribution of the groups

The “Column – Caller” is anonymously provided in the dataset and hence it is difficult to comprehend. As seen in Figure 17 Caller Data anonymously provided in the dataset

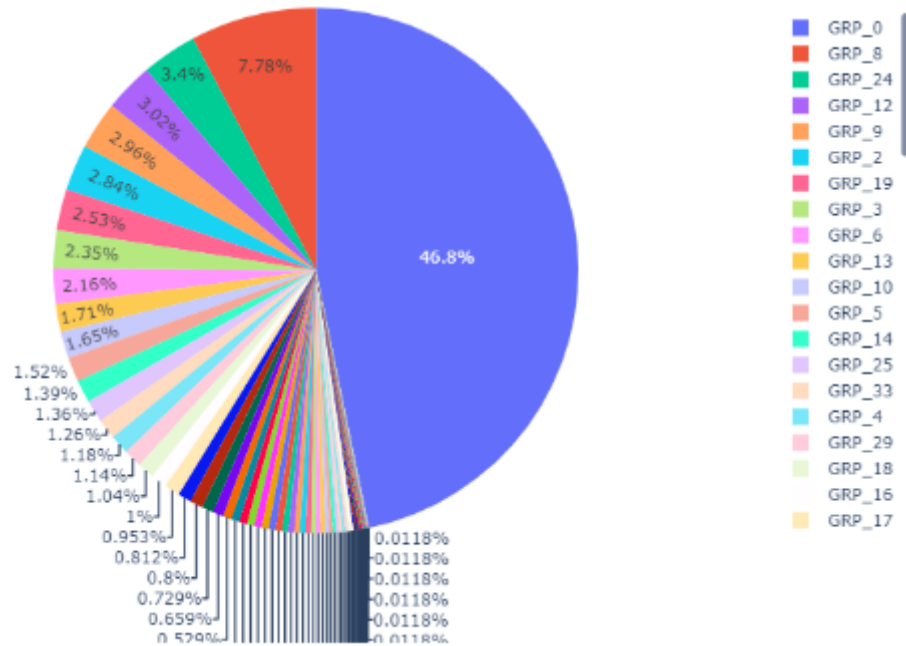


Figure 16: pie chart representation of Groups with Threshold VALUE = 50

The Caller data is then renamed with the word 'Caller' followed by an incremental number depending on the frequency of assignment. Eg: Caller1, Caller2, etc., Caller 1 has maximum ticket logged. The top 10 Caller information is displayed as below.

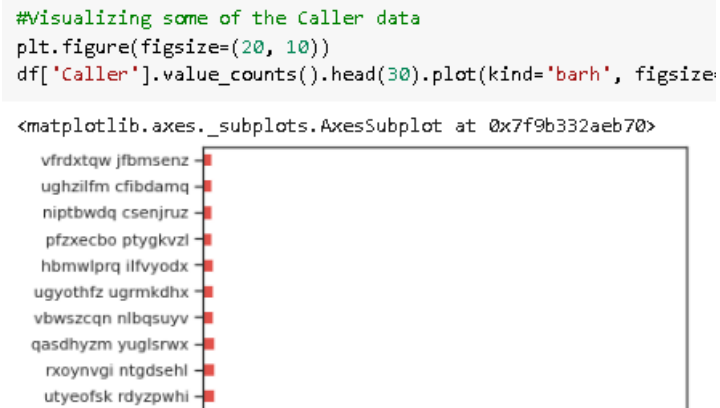


Figure 17: Caller Data anonymously provided in the dataset

This graph depicts the Caller data with frequency of tickets logged.

```
#Since caller column contains anonymous data, assigning name Caller1, Caller2,
count = 0
new_caller = []
while count != len(data):
    new_caller.append('Caller'+str(count+1))
    count = count + 1
data['caller'] = new_caller
data = data.head(20)
data.head(10)
```

	caller	n_caller
0	Caller1	810
1	Caller2	151
2	Caller3	134
3	Caller4	87
4	Caller5	71
5	Caller6	64
6	Caller7	63
7	Caller8	57
8	Caller9	54
9	Caller10	51

Figure 18: Caller frequency

Analysing the caller data, we see that one caller has reported 810 tickets, and the distribution is skewed. This could be a batch job or an anomaly and should be considered for discussion with domain expert. We are not taking up further analysis of caller for this project, for now.

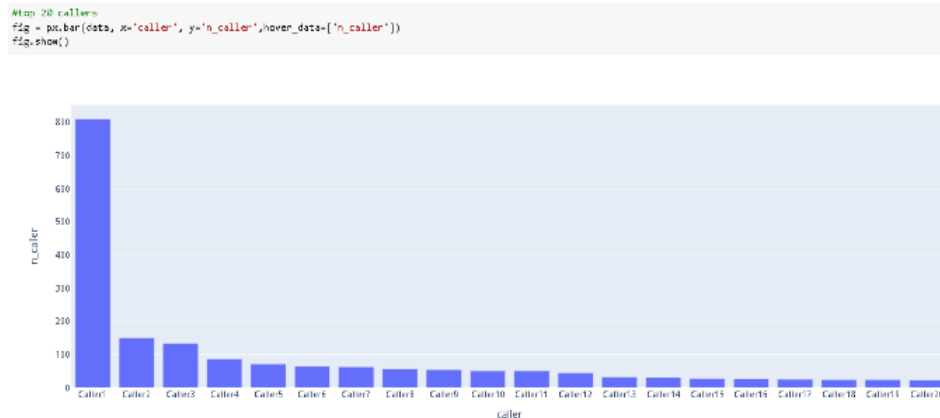


Figure 19: Barplot representing the Caller Frequency Data for Top 20 Callers

4 Data Preprocessing

The methods employed to clean our data are used from the NLTK and Regular Expression (re) library.

```
# NLTK Stop words
import nltk
nltk.download('stopwords')
nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')
nltk.download('words')
words = set(nltk.corpus.words.words())
from nltk.corpus import stopwords
stop_words = stopwords.words('english')
stop_words.extend(['received from', 'hi', 'hello', 'i', 'am', 'cc', 'sir', 'good morning', 'gentles', 'dear', 'kind', 'best', 'please', ''])
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
from gensim.utils import tokenize
```

Figure 20: Utilizing the NLTK and re Library to clean the data

4.1 Addressing Noise

STOP-WORDS

The text dataset includes many words such as ‘a’, ‘about’, ‘about’, ‘after’, ‘again’,..., which do not add any significance in model building. The most common technique to deal with these words is to remove them from the texts and documents.

PUNCTUATIONS REMOVAL

Text documents generally contains characters like punctuations or special characters and they are not necessary for text mining or classification purposes. Even though punctuation is critical to understand the meaning of the sentence, it can affect the classification algorithms negatively. Noise that is addressed as a part of this section includes – removal of brackets, newline, multi-line

spaces, numeric and alpha numeric, non-ascii text, underscores, email addresses, and disclaimers.

```
import string
import re

# Function for Text Cleaning with regex. Pass the column
def text_preprocessing(df_column):
    data = df_column.values.tolist() # Convert to list
    temp = []
    for sentence in data:
        sentence = sentence.replace("select the following link to view the disclaimer in an alternate language", '') # r
        # remove disclaimer text
        sentence = re.sub(r"\[.*?\]", " ", sentence) # remove text in []
        sentence = re.sub(r"\(.*?\)", " ", sentence) # remove text in ()
        sentence = re.sub(r"([h][t][t][p]|\S+)([w][w][w]|\S+)([\S]+[\S]+[\S]+)", " ", sentence) # remove email addresses, we
        # b address and urls
        sentence = re.sub(r"[\S]+[\d]+[\S]+", " ", sentence) # remove alphanumerics and numerics (dates, time, request id
        # etc.)
        sentence = re.sub(r"\W(?!['. ])", " ", sentence) # remove all non words with negative look back except ('. space
        # s)
        sentence = re.sub(r"[^a-zA-Z. ]+", " ", sentence) # remove non-alphabetic text
        sentence = re.sub(r"[_]", " ", sentence) # remove underscores
        sentence = re.sub(r"[\S]+", " ", sentence) # replace multiple spaces with single space
        sentence = sentence.strip('\n')
        sentence = sentence.lower()
        temp.append(sentence)
    return(temp)
```

Figure 21: Noise Removal from text

CAPITALIZATION

Sentences can contain a mixture of uppercase and lowercase words. In our dataset, we do not want to create different tokens for capitalized words. Hence we change the words to follow lowercase this brings all words in a document in same space. In some cases where capitalization changes the meaning of some words, such as "US" to "us" where first one represents the United States of America and second one is a pronoun, named identity, slang or abbreviation converters can be applied.

CONCATENATION OF ALL TEXT TO A NEW FIELD

'Summary'. The preprocessed data is then concatenated into a single column called 'Summary'.

LANGUAGE TRANSLATIONS

It has been observed that languages besides English are present in the dataset. As a part of the text processing activity, English alone has been considered and any other non-English text is dropped. However, the translation will be addressed as a part of Milestone -2.

4.2 Create Word Vocabulary Tokens

TOKENIZATION

Tokenization refers to text segmentation or lexical analysis where the large chunk of text or sentence is split into words, phrases, symbols, or elements called tokens. The main goal of this step is to extract individual words in a sentence. Along with text classification, in text mining, it is necessary to incorporate a parser in the pipeline which performs the tokenization of the documents; for example:
sentence: cant log into vpn
tokens: 'cant', 'log', 'into', 'vpn'

LEMMATIZATION

Text lemmatization is the process of eliminating redundant prefix or suffix of a word and extracting


```

1 if DELETE_CALLER:
2     df_v1["Summary"] = df_v1['Short_description'].str.cat(df_v1['Description'], sep = ". ")
3     df_v1["Summary"] = df_v2['Summary'].str.cat(df_v1['Caller'], sep = ". ")
4 else:
5     df_v1["Summary"] = df_v1['Short_description'].str.cat(df_v1['Description'], sep = ". ")
6
7 df_v2 = df_v1[['Group', 'Summary']]
8 df_v2.head(5)

```

	Group	Summary
0	GRP_0	login issue. verified user details. checked t...
1	GRP_0	outlook. received from hello team my meetings...
2	GRP_0	cant log in to vpn. received from hi i cannot...
3	GRP_0	unable to access hr tool page. unable to acces...
4	GRP_0	skype error . skype error

ERROR! Session/line number was not unique in database. History logging moved to new session 71

Figure 22: Concatenation of text into a new field – ‘Summary’

```

3 # Remove stopwords
4 df_v2['Summary'] = df_v2['Summary'].apply(lambda x: ' '.join([word for word in x.split() if word not in (stop_words)]))
5
6 # Remove words not in English Dictionary (typos, anonymised names)
7 df_v2['Summary'] = df_v2['Summary'].apply(lambda x: ' '.join([word for word in x.split() if word in (words)]))
8
9 # Tokenise 'Summary' column
10 data = df_v2.Summary.values.tolist()
11
12 data = [list(tokenize(sentences)) for sentences in data]
13
14 token_data = data
15

```

Figure 23: Utilizing Tokenize to create word tokens

the base word (lemma).

Eg: word: see or saw

lemma text: see or saw depending on whether the use of the token was as a verb or a noun

```
# lemmetise words
wordnet_lemmatizer = WordNetLemmatizer()
temp = []
for eachrow in data:
    lemma_words = []
    for eachword in eachrow:
        eachword = wordnet_lemmatizer.lemmatize(eachword, pos = "n")
        eachword = wordnet_lemmatizer.lemmatize(eachword, pos = "v")
        eachword = wordnet_lemmatizer.lemmatize(eachword, pos = ("a"))
        lemma_words.append(eachword)
    temp.append(lemma_words)
```

Figure 24: Lemmatizing words

WEIGHTED WORDS

The most basic form of weighted word feature extraction is TF (Term Frequency), where each word is mapped to a number corresponding to the number of occurrences of that word in the whole corpora. Methods that extend the results of TF generally use word frequency as a Boolean or logarithmic scaled weighting.

In all weight words methods, each document is translated to a vector (with length equal to that of the document) containing the frequency of the words in that document. Although this approach is intuitive, it is limited by the fact that words that are commonly used in the language may dominate such representations.

TERM FREQUENCY-INVERSE DOCUMENT FREQUENCY (TF-IDF)

Different word embedding procedures have been proposed to translate these unigrams into consumable input for machine learning algorithms. A very simple way to perform such embedding is weighted words term-frequency (TF) and TF-IDF where each word will be mapped to a number corresponding to the number of occurrence of that word in the whole corpora.

Although tf-idf tries to overcome the problem of common terms in document, it still suffers from some other descriptive limitations. Namely, tf-idf cannot account for the similarity between words in the document since each word is presented as an index.

```
# Create Weighted Word Vectors
tfidf_vectors = TfidfVectorizer(min_df=3,max_features= maxlen)
tfidf_db = tfidf_vectors.fit_transform(data).toarray()
tfidf_db = pd.DataFrame(tfidf_db)

tfidf_db.head()
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24
0	0.214291	0.000000	0.0	0.0	0.0	0.0	0.0	0.286802	0.0	0.0	0.0	0.278404	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.000000	0.000000	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.530371	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	0.000000	0.000000	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	0.000000	0.54983	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.000000	0.000000	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.000000	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

5 rows x 25 columns

Figure 25: Creating weighted vectors of the vocabulary

5 Train Test Split and Evaluation

Before we proceed to build the model, we encode the Group information labelled as GRP_X as 0,1... using label encoder(). This field represents the Target column.

```
le = preprocessing.LabelEncoder()
df_v2['Group'] = le.fit_transform(df_v2['Group']) # LabelEncode 'Groups'
df_v2.head(20)
```

Figure 26: Encoding the Group Data using LabelEncoder

TRAIN TEST SPLIT

We then map the X(input data) to the vectorized data and y(target) to encoded groups in order to build our models. The training and testing datasets are formulated from the X and y data, in a ratio of 60-40 training and test data respectively using the train_test_split() function. The ideal range is 60-40 to 80-20. On changing the test_size parameter, we can modify the testing data size.

```
X = tfidf_db
y = df_v2['Group']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.40, random_state=42)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=42)
```

Figure 27: Setting X and y(Target), splitting between training and testing sets

We then fit the training data to the model using model.fit(X_train, y_train). The predicted value of y is obtained using model.predict(X_test).

DECIDING EVALUATION CRITERIA

Training and Testing scores are determined based on the training and testing sets respectively for the models that are build. We use `model.score(X,y)` to evaluate the same. Recall, precision, and confusion matrix are metrics used to determine the algorithms response to the multi-classification problem. Recall (also known as sensitivity) is the fraction of positives events that you predicted correctly as shown below and Precision is the fraction of predicted positives events that are actually positive as shown below. Confusion matrix is the matrix to the right which depicts the Y_actual and Y_predicted.

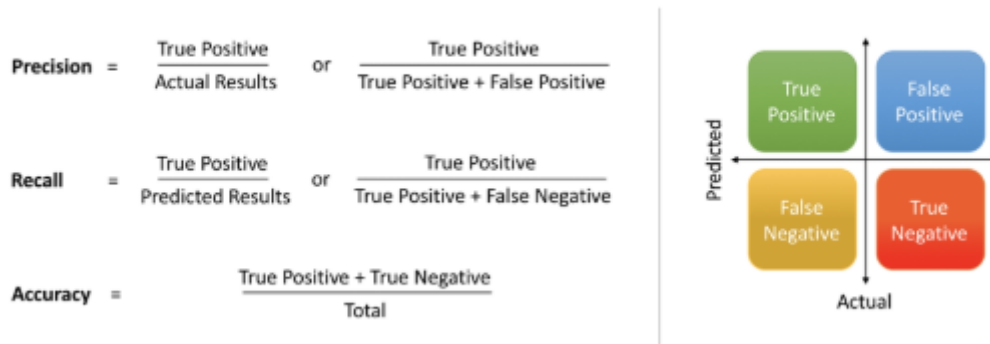


Figure 28: Recall, precision, and confusion matrix

TRADE-OFF

F1 score is the harmonic mean of recall and precision, with a higher score as a better model. The F1 score is calculated using the following formula:

$$\text{F1 Score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

We have employed weighted-average F1-score, or weighted-F1 wherein we weight the F1-score of each class by the number of samples from that class.

6 Deciding Models and Model Building

The following models are considered as a part of classifying the tickets into their respective groups. In Milestone 1, we employed Traditional classification algorithms, such as the Naïve Bayes, Support Vector Classifier, Decision Trees, Random Forests, and Ensemble. For the given problem which is a supervised learning, multi-classification problem, the approach was to tackle the problem using conventional techniques in Milestone 1 and thereby proceeding to Deep Neural networks, such as -Long Short-Term Memory (LSTM), Recurrent Convolutional Neural Networks (RCNN), Random Multimodel Deep Learning (RMDL), etc in Milestone 2.

Majority of the time (about 60% was spent in Data cleaning activities to prepare the data for modelling.

6.1 Naive Bayes

Naive Bayes is a probabilistic learning algorithm derived from Bayes Theorem. Naive Bayes Model is considered to be extremely fast, reliable, and has stable classification ability relative to other classification algorithms. The algorithm is based on the assumption that each feature is independent of each other while predicting the classification.

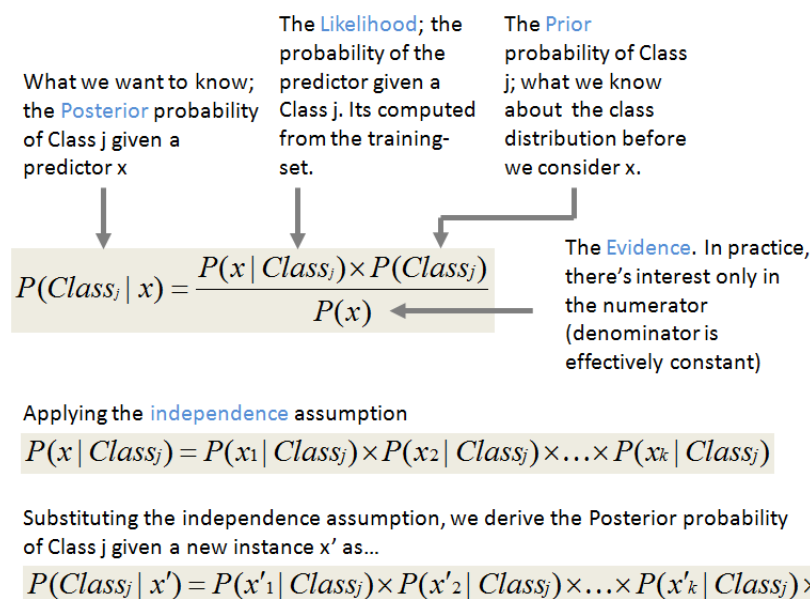


Figure 29: Explanation of Naive Bayes, Source: <http://shatterline.com/blog/2013/09/12/not-so-naive-classification-with-the-naive-bayes-classifier/>

PROS:

- Simple, fast and well in multi class prediction.
- Performs better with less training data as it assumes feature independence

CONS:

- Bad estimator hence the probability outputs are not taken too seriously.
- Assumptions of independent feature cannot represent real time data.
- Zero frequency - If training data set gets a category not trained on earlier, then model will assign a 0 (zero) probability and will be unable to make a prediction.

The multinomial Naive Bayes classifier is suitable for classification with discrete features (e.g., word counts for text classification), we will build Multinomial Naive Bayes model for our dataset.

Multinomial Naive Bayes

```
[42] 1 # Naive Bayes
      2 NBModel = MultinomialNB(alpha = 0.001)
      3 NBModel.fit(X_train, y_train)
      4 NB_y_pred = NBModel.predict(X_test)
      5 print('NB Training Accuracy:', 100*NBModel.score(X_train , y_train))
      6 print('NB Test Accuracy:', 100*NBModel.score(X_test , y_test))

NB Training Accuracy: 57.38672286617492
NB Test Accuracy: 55.34134007585335

[43] 1 print ('Precision Score:', precision_score(y_test, NB_y_pred, pos_label='Positive', average='weighted'))
      2 print ('Recall Score:', recall_score(y_test, NB_y_pred, pos_label='Positive', average='weighted'))
      3 F1score = f1_score(NB_y_pred, y_test, average='weighted')
      4 print('F1 Score:', F1score)
```

Figure 30: Implementation of Naive Bayes

RESULT: We can see that the Training accuracy is 57% and testing accuracy is 55% with Naive Bayes Model. The model is able to predict True Positives and False Negatives equally.

6.2 Support Vector Classifier (SVC)

Support Vector Machine (SVM) creates a hyperplane between the classes which acts as decision boundary for each class. Data falling within these boundaries will belong to that particular class. SVM can classify non-linear data and can capture complex relationships between data points without having to perform difficult transformations. While Naïve Bayes treats the features of dataset as independent, SVM analyses the interactions between each feature to certain.

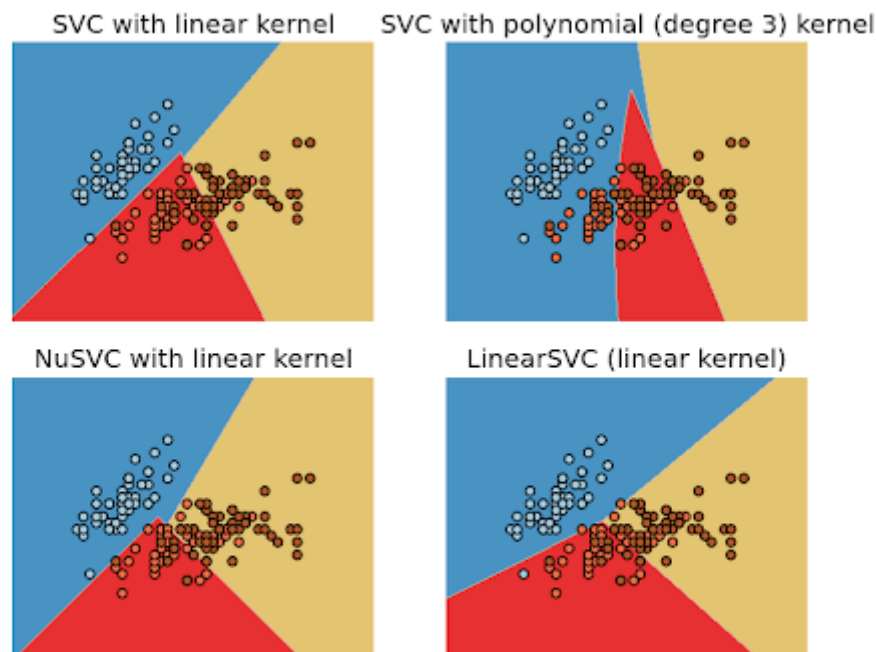


Figure 31: Visual representation of SVC Source: <http://scikit-learn.sourceforge.net/0.8/modules/svm.html>

We use linear kernel in SVC for this problem as the text are transformed into higher dimensions by using TF IDF and the data in higher dimensions are linearly separable. The linear kernel equation for predicting a new input using the dot product between the input (x) and each support vector (x_i) is calculated as $f(x) = B(0) + \sum(a_i * (x, x_i))$

PROS:

- Less affected by outliers, relatively computationally efficient and accurate than its competitors.
- Effective where number of features are greater than the number of samples.

- Good generalization capabilities which prevents it from over-fitting

CONS:

- Does not perform very well when the data of target classes are overlapping
- Choosing an appropriate Kernel function for handling the non-linear data could be tricky and complex
- Requires lot of memory size to store all support vectors and takes long time to train on larger dataset

▼ Support Vector Machine



```

1 # Creating SVC Model
2 svm_model = SVC(kernel='linear',C=10)
3 svm_model.fit(X_train, y_train)
4 y_pred = svm_model.predict(X_test)
5 print('Training Accuracy:', 100*svm_model.score(X_train , y_train))
6 print('Test Accuracy:',100*svm_model.score(X_test , y_test))
7 |

Training Accuracy: 69.06217070600633
Test Accuracy: 55.56257901390644

[ ] 1 print ('Precision Score:', precision_score(y_test, y_pred, pos_label='Positive', average='weighted'))
2 print ('Recall Score:', recall_score(y_test, y_pred, pos_label='Positive', average='weighted'))
3
4 F1score = f1_score(y_pred, y_test, average='weighted')
5 print('F1 Score:', F1score)
6 print(confusion_matrix(y_test,y_pred))
7 print(classification_report(y_test,y_pred))

```

Figure 32: Implementation of SVC

RESULT:

We can see that the Training accuracy is around 69% and testing accuracy is around 55% with SVC. The model is able to predict True Positives and False Negatives equally.

6.3 Decision Tree

Decision Tree solves the problem of machine learning by transforming the data into tree representation. Each internal node of the tree denotes an attribute and each leaf node denotes a class label. Decision tree algorithm can be used to solve both regression and classification problems. Decision Tree creates a training model which can use to predict class by learning decision rules inferred from training data. Decision tree creates a model to predict the labels by learning the decision rule from training data.

The cost functions try to find most homogeneous branches, or branches having groups with similar responses. The mean of responses of the training data inputs of that group is considered as prediction for that group.

Classification: $G = \sum(pk * (1 - pk))$

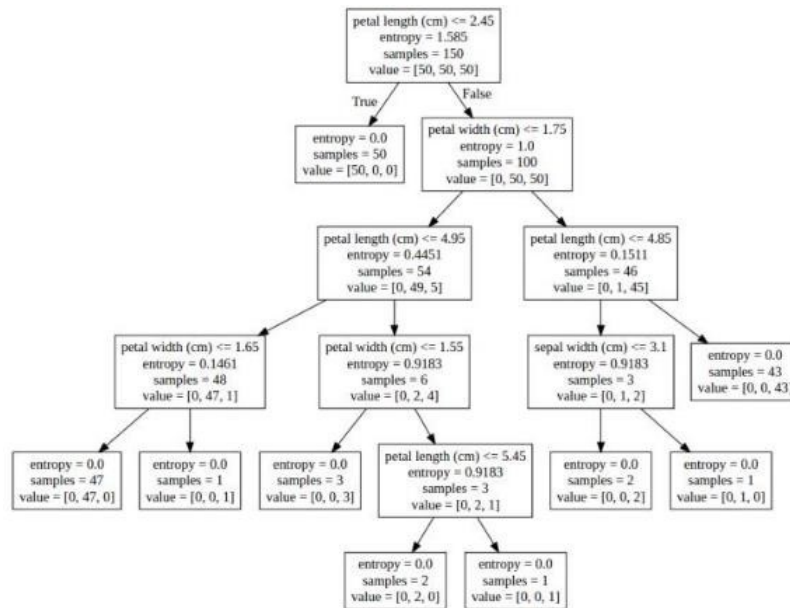


Figure 33: Explanation of Decision Trees source: <https://www.kdnuggets.com/2017/05/simplifying-decision-tree-interpretation-decision-rules-python.html>

PROS:

- Missing values does not affect decision tree.
- Requires less effort for data preparation during pre-processing, does not need scaling or normalization.
- The Number of hyper-parameters to be tuned is almost null.
- A Decision trees model is very intuitive and Interpretation of a complex Decision Tree model can be simplified by its visualizations

CONS:

- Instable as a small change in the data can cause a large change in the structure of the decision tree.
- For a Decision tree sometimes calculation can go far more complex compared to other algorithms.
- Probability of overfitting is high and training time is more making it expensive and complex.

Decision Tree

```
[46] 1 # Using Decision Tree Classifier
      2 dt_model = DecisionTreeClassifier(criterion = 'entropy' )
      3 dt_model.fit(X_train, y_train)

DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='entropy',
                      max_depth=None, max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')

[47] 1 y_predict = dt_model.predict(X_test)
      2 print('Training Accuracy:',100*dt_model.score(X_train, y_train))
      3 print('Test Accuracy:',100*dt_model.score(X_test , y_test))

Training Accuracy: 83.60379346680716
Test Accuracy: 50.50568900126422
```

Figure 34: Implementation of Decision Trees

- low prediction accuracy compared to other algorithms and can become complex when there are many class labels.

RESULT:

We can see that the training accuracy is around 83% and testing accuracy is around 50% with Decision Trees. The model is able to predict True Positives and False Negatives almost equally.

6.4 Random Forest (RF)

Random forest classifier creates a number of decision trees from randomly selected subset of training set. It then uses averaging to improve the predictive accuracy and control over-fitting. Random forest applies weight concept, tree with high error rate are given low weight value and vice versa. This would increase the decision impact of trees with low error rate.

PROS:

- Random forest is an accurate and robust method because of the number of decision trees participating in the process.
- It takes the average of all the predictions, which cancels out the biases thereby does not suffer from the overfitting problem.
- Can handle missing values and can be used in both classification and regression problems.

CONS:

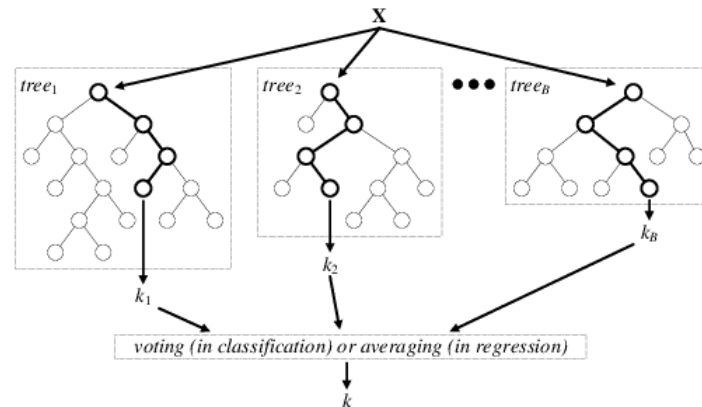


Figure 35: Explanation of Random Forests (Source: <https://www.researchgate.net>)

- Random forest is slow in generating predictions as it has multiple decision trees. All the trees in the forest must make a prediction then perform voting on it. This whole process is time-consuming.
- The model is difficult to interpret compared to a decision tree.

Random Forest Classifier

```
[48] 1 #Using Random Forest Classifier
      2 rfcl = RandomForestClassifier(criterion = 'entropy', n_estimators = 50)
      3 rfcl = rfcl.fit(X_train, y_train)
      4 test_pred = rfcl.predict(X_test)
      5 print('Training Accuracy:', 100*rfcl.score(X_train, y_train))
      6 print('Test Accuracy:', 100*rfcl.score(X_test, y_test))

Training Accuracy: 83.60379346680716
Test Accuracy: 55.94184576485461

[49] 1 print ('Precision Score:', precision_score(y_test, test_pred, pos_label='Positive', average='weighted'))
      2 print ('Recall Score:', recall_score(y_test, test_pred, pos_label='Positive', average='weighted'))
      3
      4 F1score = f1_score(test_pred, y_test, average='weighted')
      5 print('F1 Score:', F1score)
      6 print(confusion_matrix(y_test, test_pred))
      7 print(classification_report(y_test, test_pred))
```

Figure 36: Implementation of Random Forests

RESULT:

We can see that the training accuracy is around 83% and testing accuracy is around 55% with Random Forests. The model is able to predict True Positives and False Negatives almost equally.

6.5 Ensemble

Most of the errors from a model's learning are from three main factors: variance, noise, and bias. Using ensemble methods, we can increase the stability of the final model and reduce the errors mentioned previously. By combining many models, we can (mostly) reduce the variance, even when they are individually not great, as we will not be affected by random errors from a single source. We have used two common ensemble techniques to solve the Automated Ticket Classification problem assigned in this project

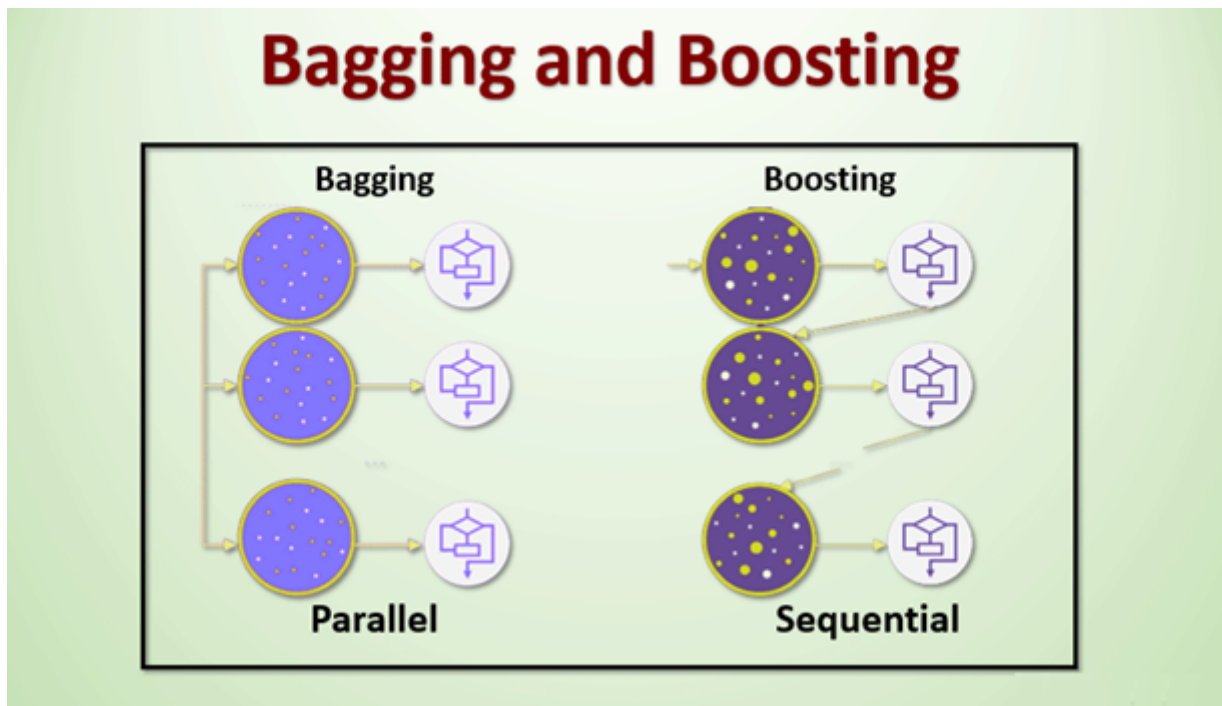


Figure 37: Bagging and Boosting Source: <https://www.educba.com/bagging-and-boosting/>

BAGGING

Bagging is shorthand for the combination of bootstrapping and aggregating. Bootstrapping is a method to help decrease the variance of the classifier and reduce overfitting by resampling data from the training set with the same cardinality as the original set. In bagging there is a tradeoff between base model accuracy and the gain you get through bagging. The aggregation from bagging may improve the ensemble greatly if there is an unstable model. Once the bagging is done, and all the models have been created on (mostly) different data, a weighted average is then used to determine the final score.

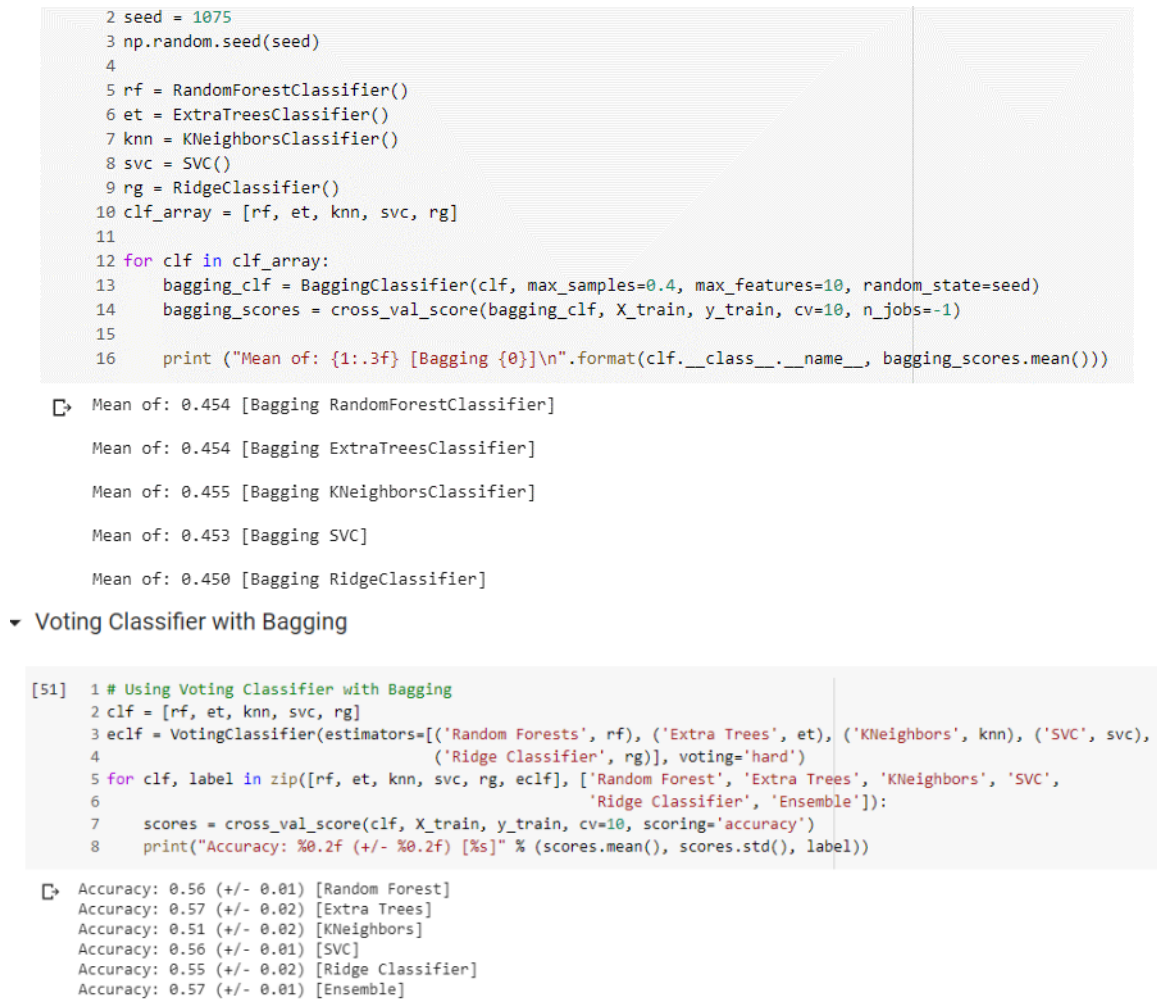


Figure 38: Ensemble model -Bagging

BOOSTING

The main idea of boosting is to add additional models to the overall ensemble model sequentially. With each iteration of boosting, a new model is created and the new base-learner model is trained (updated) from the errors of the previous learners. The algorithm creates multiple weak models whose output is added together to get an overall prediction and the boosted gradient shifts the current prediction nudging it to the true target. The gradient descent optimization occurs on the output of the various models, and not on their individual parameters.

For Bagging we have used an ensemble of RandomForestClassifier, ExtraTreeClassifier, KNeighborsClassifier, SVM and RidgeClassifier with hard voting, which just need a majority of classifiers to determine what the result could be. Please see Figure 38 below. We can see that the training accuracy is around 84% and testing accuracy is around 58% with Random Forests. The model is able to predict True Positives and False Negatives almost equally. We can see that the training

accuracy is around 84% and testing accuracy is around 58% with Random Forests. The model is able to predict True Positives and False Negatives almost equally.

RESULT:

We can see that with Bagging ensemble model, the mean of each classifier used is around 48% and the maximum accuracy when using hard voting is 59% with the chosen classifiers. For Boosting we have used an ensemble of AdaBoostClassifier, GradientBoostingClassifier and XGBClassifier using hard voting from EnsembleVoteClassifier(mlxtend).

```
Boosting Classifier

[ ] 1 # Using Boosting Classifiers
    2
    3 ada_boost = AdaBoostClassifier()
    4 grad_boost = GradientBoostingClassifier()
    5 xgb_boost = XGBClassifier()
    6
    7 boost_array = [ada_boost, grad_boost, xgb_boost]
    8 eclf = EnsembleVoteClassifier(clfs=[ada_boost, grad_boost, xgb_boost], voting='hard')
    9 labels = ['Ada Boost', 'Grad Boost', 'XG Boost', 'Ensemble']
   10 for clf, label in zip([ada_boost, grad_boost, xgb_boost, eclf], labels):
   11     scores = cross_val_score(clf, X_train, y_train, cv=10, scoring='accuracy')
   12     print("Mean: {:.3f} of [{}]" .format(scores.mean(), label))

☞ Mean: 0.486 of [Ada Boost]
   Mean: 0.539 of [Grad Boost]
   Mean: 0.560 of [XG Boost]
   Mean: 0.556 of [Ensemble]
```

Figure 39: Ensemble model – Boosting

RESULT:

The maximum accuracy is with 58% with XGBoost classifier. Out of the ensemble models built, bagging with hard voting seem to give slightly higher accuracy of 59

7 Improving Model Performance

Hyperparameters in Machine Learning are user-controlled settings of ML models. These hyperparameters influence how the model's parameters are updated and learned during the training. If the right hyperparameters are set, the model will learn the most optimal weights that it can with a given training algorithm and data. The best hyperparameters are found by trial and error method manually or with Grid Search module

7.1 Naive Bayes

```
[ ] 1 def find_optimal_k(X,y, lr_list):
    2
    3     # empty list that will hold cv scores
    4     scores = []
    5     for lr in lr_list:
    6         model_nb = MultinomialNB(alpha = lr)
    7         model = model_nb.fit(X_train, y_train)
    8         # predict the response on the crossvalidation train
    9         predict = model.predict(X_test)
   10
   11     # evaluate accuracy
   12     accuracy = accuracy_score(y_test, predict, normalize=True)
   13     scores.append(accuracy)
   14
   15     # changing to misclassification error
   16     mean_square_error = [1 - x for x in scores]
   17
   18     # determining best alpha
   19     optimal_alpha = lr_list[mean_square_error.index(min(mean_square_error))]
   20     print('\nThe optimal alpha is ', optimal_alpha)
```

Figure 40: Hyperparameter tuning – Naïve Bayes

RESULT:

The best alpha found is 0.00001. The best training accuracy is 57% and test accuracy is 55%. With shufflesift, we found that the least fit time is 0.021 seconds

7.2 Support Vector Classifier

▼ Support Vector Machine

The kernel used for SVM in this project is 'linear'. We can use grid search to get optimal parameters that gives best accuracy

```
[ ] 1 #Set parameters for grid search
    2 param = {
    3   'C': (np.arange(1.4,2.0,0.2)) , 'kernel': ['linear'],
    4   'C': (np.arange(1.4,2.0,0.2)) , 'gamma': [0.01,0.03,0.05], 'kernel': ['rbf'],
    5   'degree': [2,3,4] , 'gamma':[0.01,0.03,0.05], 'C':(np.arange(0.1,1,0.1)) , 'kernel':['poly']
    6   }
    7
    8 #import Grid Search
    9 from sklearn.model_selection import GridSearchCV
   10
   11 model_svm = GridSearchCV(svm_model, param, cv=10, scoring='accuracy')

[ ] 1 model_svm.fit(X_train, y_train)
    2 print(model_svm.best_score_)
    3 print(model_svm.best_params_)

0.4322447257383966
{'C': 0.1, 'degree': 2, 'gamma': 0.01, 'kernel': 'poly'}
```

Figure 41: Grid Search - SVM

RESULT:

The best param found is C-0.1, degree-2, gamma-0.01, kernel-poly. The best accuracy is 43

7.3 Decision Tree

```
Decision Tree

[64] 1 #setting parameters for grid search
      2 max_dep_range = [4,5,8,10, 12,15,18,20,25,30, 40,50,60,70,80]
      3 min_lf = np.arange(4,20,2)
      4 tree_param = [{'criterion': ['entropy', 'gini'], 'max_depth': max_dep_range}, {'min_samples_leaf': min_lf}]
      5 GD = GridSearchCV(dt_model, tree_param)
      6 GD.fit(X_train,y_train)
      7 print("Best Hyper Parameters for DT:\n",GD.best_params_)

Best Hyper Parameters for DT:
{'criterion': 'entropy', 'max_depth': 15}

[65] 1 #building model with best parameters
      2 DT_model = DecisionTreeClassifier(criterion = 'entropy', max_depth= 15)
      3 DT_model.fit(X_train,y_train)
      4 test_pred_DT = DT_model.predict(X_test)
      5 print('DT train accuracy after Grid Search:', DT_model.score(X_train , y_train))
      6 print('DT test accuracy after Grid Search:', DT_model.score(X_test , y_test))

DT train accuracy after Grid Search: 0.6112871287128713
DT test accuracy after Grid Search: 0.5518265518265518
```

Figure 42: Decision Tree – Grid Search

RESULT:

The training accuracy is 61% and testing accuracy is 55% after running the model with Grid search's best parameters. The train accuracy has come down while the test accuracy is slightly more than the initial model. The model seems to be able to generalize better but the test accuracy is still not good enough

7.4 Random Forest

Using Grid Search for Random Forest model, we see that the best hyperparameters for the model is Entropy and 100 estimators. Applying these parameters and running the model gives us test accuracy of 58% but the accuracy is not a great improvement from the initial model accuracy.

Random Forest

```
[60] 1 param_grid = {'criterion':['gini','entropy'],'n_estimators':[50,100,150,200,250]}
      2 gs = GridSearchCV(rfcl,param_grid)
      3 gs
```

```
1 #get best parameter
2 gs.fit(X_train,y_train)
3 print("Best Hyper Parameters:\n",gs.best_params_)
```

```
Best Hyper Parameters:
{'criterion': 'entropy', 'n_estimators': 100}
```

```
[62] 1 #Build Modl with best parameters got from Grind Search
      2 rfcl = RandomForestClassifier(criterion = 'gini', n_estimators = 150)
      3 rfcl = rfcl.fit(X_train, y_train)
      4 test_pred = rfcl.predict(X_test)
      5 acc_RFGS = accuracy_score(y_test, test_pred)
      6 print('Training Accuracy:', 100*rfcl.score(X_train , y_train))
      7 print('Test Accuracy:',100*rfcl.score(X_test , y_test))
```

```
Training Accuracy: 84.73267326732673
Test Accuracy: 58.83575883575883
```

Figure 43: Random Forest – Grid Search

8 Deep Learning Networks

In the basic models built so far the maximum accuracy we could reach was 58% with hyper tuning Random Forest model. We can now check to see if we can get better result with some deep neural network models.

8.1 LSTM

Long short-term memory (LSTM) is an artificial recurrent neural network (RNN) architecture. Recurrent Neural Networks (RNN) has an internal state that store information of past inputs for a certain time by which the model learns the context of the information. Where the inputs of the model are sequential in nature, the model can return the output as sequence with the contextual information. In our project, we use the bidirectional LSTM which presents each training sequence forwards and backwards to two separate recurrent nets, both of which are connected to the same output layer. This means that for every point in a sequence, the model has complete information.

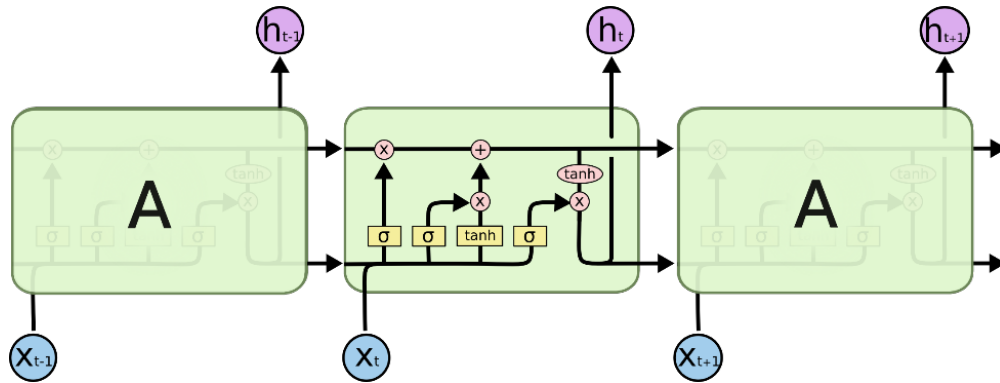


Figure 44: LSTM cell Architecture(<https://colah.github.io/posts/>)

```
[ ] 1 model = tf.keras.models.Sequential([
2     tf.keras.layers.Embedding(vocabulary_size, output_dim=100, input_length=maxlen),
3     tf.keras.layers.LSTM(128, return_sequences=True, dropout=0.2),
4     tf.keras.layers.BatchNormalization(),
5     tf.keras.layers.TimeDistributed(Dense(32, activation='tanh')),
6     tf.keras.layers.Dropout(0.25),
7     tf.keras.layers.BatchNormalization(),
8     tf.keras.layers.TimeDistributed(Dense(32, activation='tanh')),
9     tf.keras.layers.Flatten(),
10    tf.keras.layers.Dense(24, activation='softmax')
11 ])
12
13 optimiser = keras.optimizers.Adam(learning_rate=0.001, beta_1=0.9, beta_2=0.999, amsgrad=False)
14 model.compile(optimizer=optimiser, loss=tf.keras.losses.SparseCategoricalCrossentropy(), metrics=['accuracy'])

[ ] 1 model.load_weights("LSTM_Weights.best.hdf5")
2 optimiser = keras.optimizers.Adam(learning_rate=0.001, beta_1=0.9, beta_2=0.999, amsgrad=False)
3 model.compile(optimizer=optimiser, loss=tf.keras.losses.SparseCategoricalCrossentropy(), metrics=['accuracy'])
```

Figure 45: LSTM Model

9 Challenges, Approach and Mitigation

10 Code and Deliverable

- Final Report_{Group3NLP.tex}PDF format – FinalReport_{Group3NLP.pdf}
- Filenames listed, attachments in the Great Learning (GL) portal
- NLP_{AutomatedTicketAssignmentFinal.ipynb}NLP_{AutomatedTicketAssignmentFinal.html}
- Snapshot of Github Repo (Capstone)-
- <https://github.com/GLNLPGroup3/Capstone>

Challenges	Approach	Mitigation
1. Hardware		
Personal machines used for the project has limited in storage and processing power.	Find easy and free platforms with sufficient hardware and processing support.	The platform which was used to achieve the task was - Google Colab which provides around 15GB of Storage space on the Google Drive as well as its GPU which empowers us to train and test our models effectively.
Data is the most importance piece of the puzzle with regards to Machine Learning(ML) problems. Our observations		
2. Class Imbalance		
We can observe that the number of observations in each group is poorly distributed. There are totally 74 groups with some groups having as less as one observation.	Use sampling techniques would enable us to down sample the majority classes or/and upsample the minority classes. Group the minority classes into one group for classification	We created a group clubbing the minority classes into one group. Groups having less than 50 tickets will be categorized into one group. We also used upsampling with stratification and down-sampling of majority class for Neural Network model. This increased the accuracy and precision.
3. Noisy Data		
Data collected from the systems would be noisy with extra characters like punctuations, html texts, special characters etc.	Cleaning data is the primary task to any data modelling problem. We spend a considerable amount of time cleaning the data and preparing it for modelling	In the view of preparing the data for modelling, we must first clean the data. This has been accomplished using NLTK and RE (regular expressions) Libraries.
4. Multi-lingual data		
Besides English, we observed non-English text in the dataset.	We checked if libraries from Google Translate and other language translation modules would work for our purpose. However, the limitations exceeded the cause. As a part of the text processing activity, English text has been considered and any other non-English text was dropped.	For milestone 1, we considered only English text for the processing and would address the non-English text either through a translation or a mechanism to map the same using word mappings. The findings would be potentially shared in Milestone 2.

Challenges	Approach	Mitigation
5. Collaboration		
The team is located in different parts of the country and the course is completely online which posed a challenge in communication and collaboration.	As a team, the primary goal is to be able to share data between the team, communicate effectively and thereby work together.	We used the following platforms which proved efficient in meeting our team's expectations and goals. a. Github b. Google Drive c. Telephony, Whatsapp groups, etc.