

# **Automatic Ticket Classification**

NLP

**FINAL REPORT**  
**CAPSTONE PROJECT - GREAT LEARNING PGAIML**



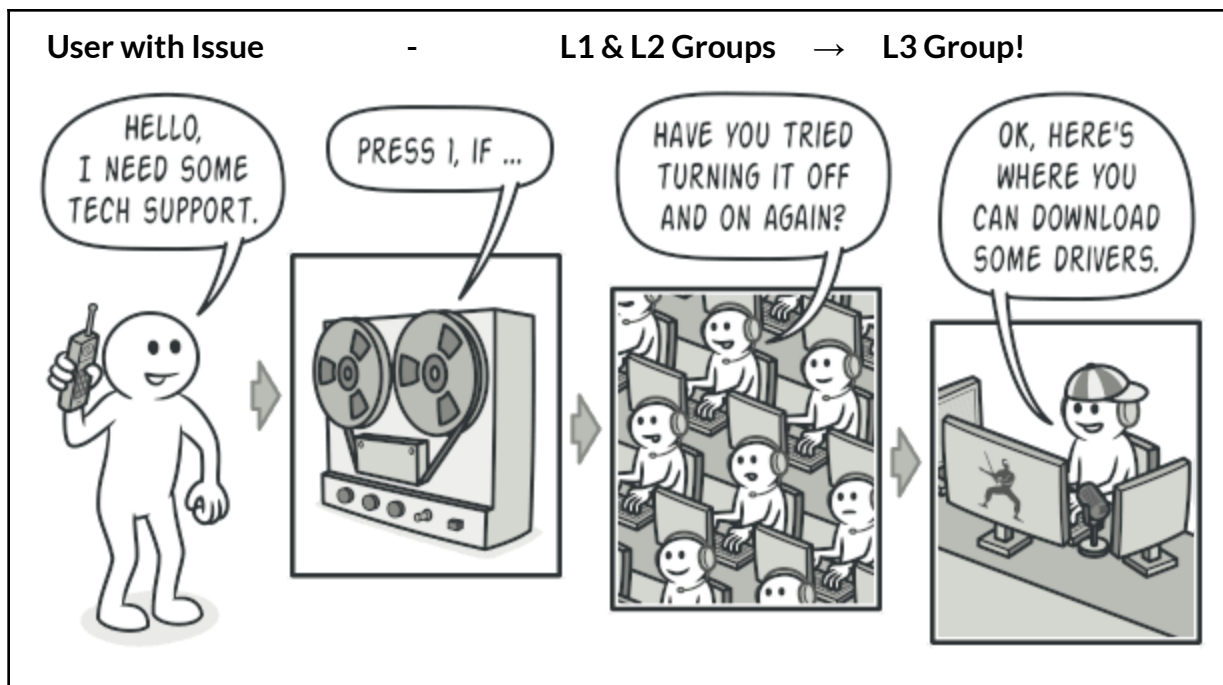
**GROUP 9**  
**MENTORED BY**  
**Vaibhav Kulkarni**

## 1. Summary of the Problem Statement, Data and Findings

### 1.1. Abstract

Excellent & Effective Customer Support is Quintessential to the running of any business organization, no matter its size—84% of organizations working to improve customer service report an increase in revenue. In the current scenario, various incidents faced by the business are all assigned to two L1/L2 teams. Only 54% of these incidents are resolved at this level. For all the rest, the incidents are escalated to L3 teams to be resolved. Additionally, the manual reassignment to various functional groups was found to have an error rate of around 25%. This added overhead cost of time and resources of re-assigning the incidents is detrimental to the customer support efficiency causing delays and bad customer experiences. A better allocation and practical usage of the functional groups' resources will result in substantial cost, time savings and better customer support overall.

Hence, we aim to build a classifier using state-of-the-art NLP techniques to classify the tickets to various functional groups by just analyzing the text of the multiple issues, thereby driving direct business value in IT customer support.



## 1.2. Dataset

- Our dataset consists of 8500 data points, each consisting of a short description of the issue, a longer description, the caller name (appears to be encrypted and anonymized in the given dataset to protect privacy) and the target class group to which the issue has to be assigned to:

```
dataset.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8500 entries, 0 to 8499
Data columns (total 4 columns):
 #   Column              Non-Null Count  Dtype  
---  -
 0   Short description    8492 non-null   object 
 1   Description          8499 non-null   object 
 2   Caller               8500 non-null   object 
 3   Assignment group    8500 non-null   object 
dtypes: object(4)
memory usage: 265.8+ KB
```

- There seem to be missing values in the Short description and Description columns, which needs to be looked into and handled. There are eight nulls/missing values present in the Short description and one null/missing value present in the description column.

```
dataset.isna().sum() # Few missing values
```

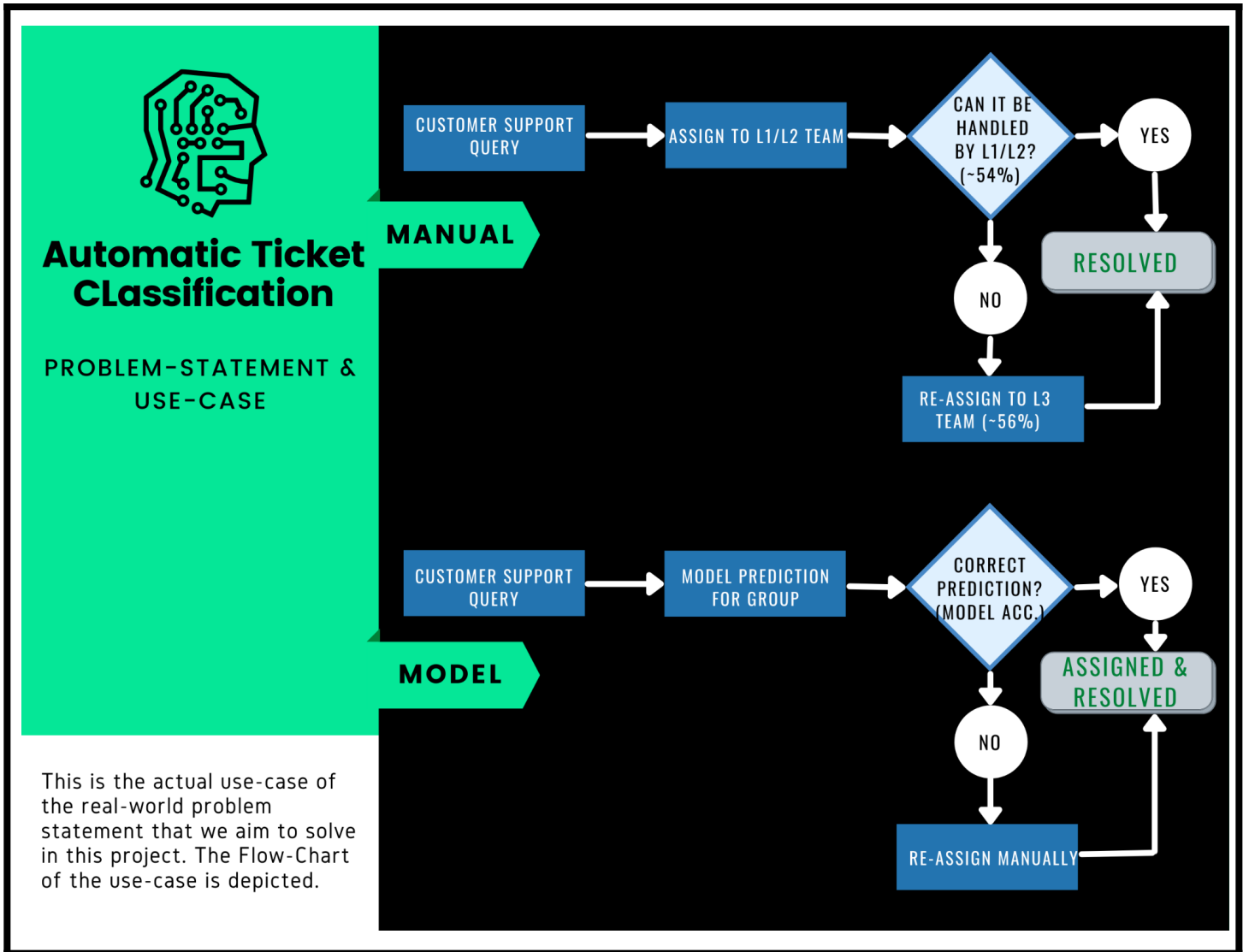
```
short_description    8
description          1
caller               0
group               0
dtype: int64
```

```
dataset[dataset.isna().any(axis=1)] # check rows with missing values
```

	short_description	description	caller	group
2604	NaN	\r\n\r\nreceived from: ohdrnswl.rezuibdt@gmail...	ohdrnswl rezuibdt	GRP_34
3383	NaN	\r\n-connected to the user system using teamvi...	qftpazns fxpnytmk	GRP_0
3906	NaN	-user unable tologin to vpn.\r\n-connected to...	awpcmsej ctdiuqwe	GRP_0
3910	NaN	-user unable tologin to vpn.\r\n-connected to...	rhwsmefo tvphyura	GRP_0
3915	NaN	-user unable tologin to vpn.\r\n-connected to...	hxripljo efzounig	GRP_0
3921	NaN	-user unable tologin to vpn.\r\n-connected to...	czadygo veiosxby	GRP_0
3924	NaN	name:vvqgbdhm fwchqjor\nlanguage:\nbrowser:mic...	vvqgbdhm fwchqjor	GRP_0
4341	NaN	\r\n\r\nreceived from: eqmunioy.ehxkcbgj@gmail...	eqmunioy ehxkcbgj	GRP_0
4395	i am locked out of skype	NaN	viyglzfo ajtfzpkb	GRP_0

- The independent features are short descriptions and descriptions, and the target/dependent feature is the group.

### 1.3. Use-Case

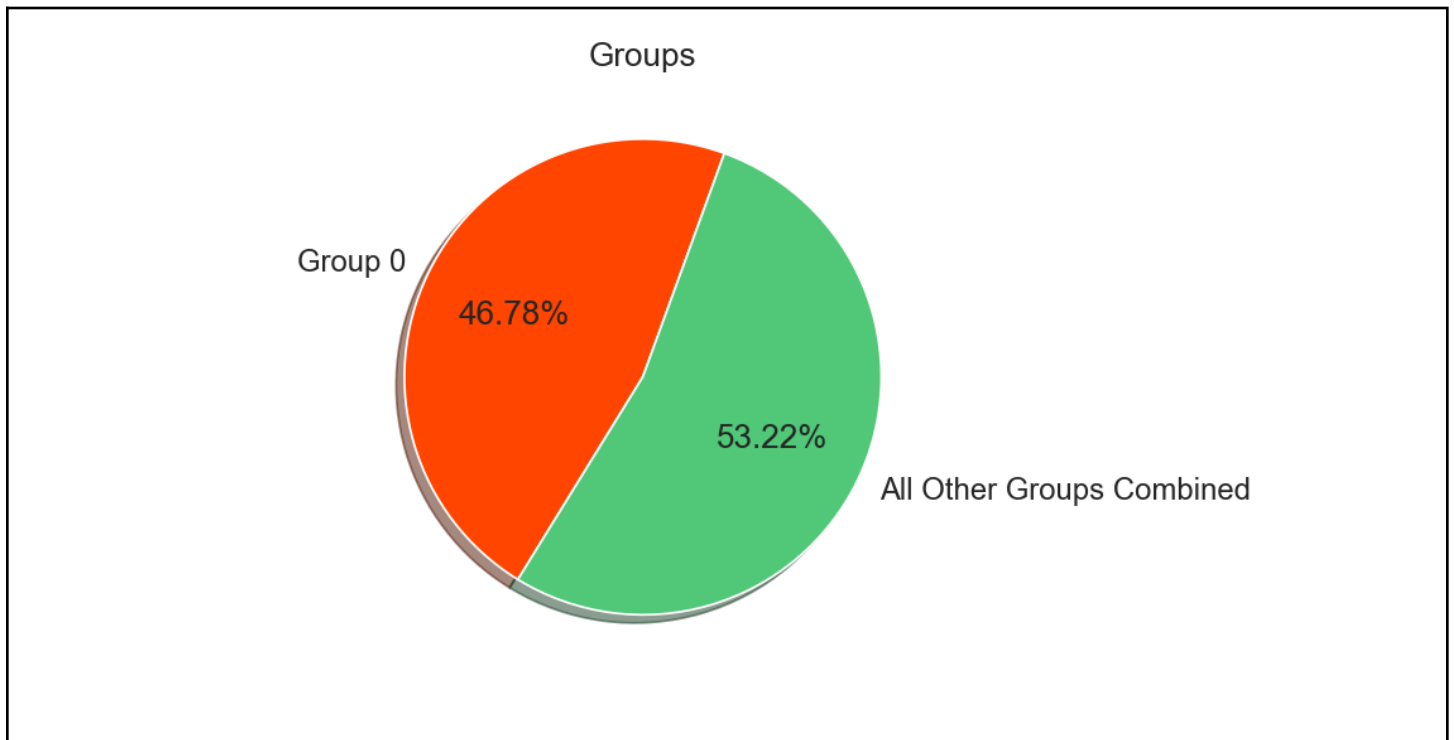


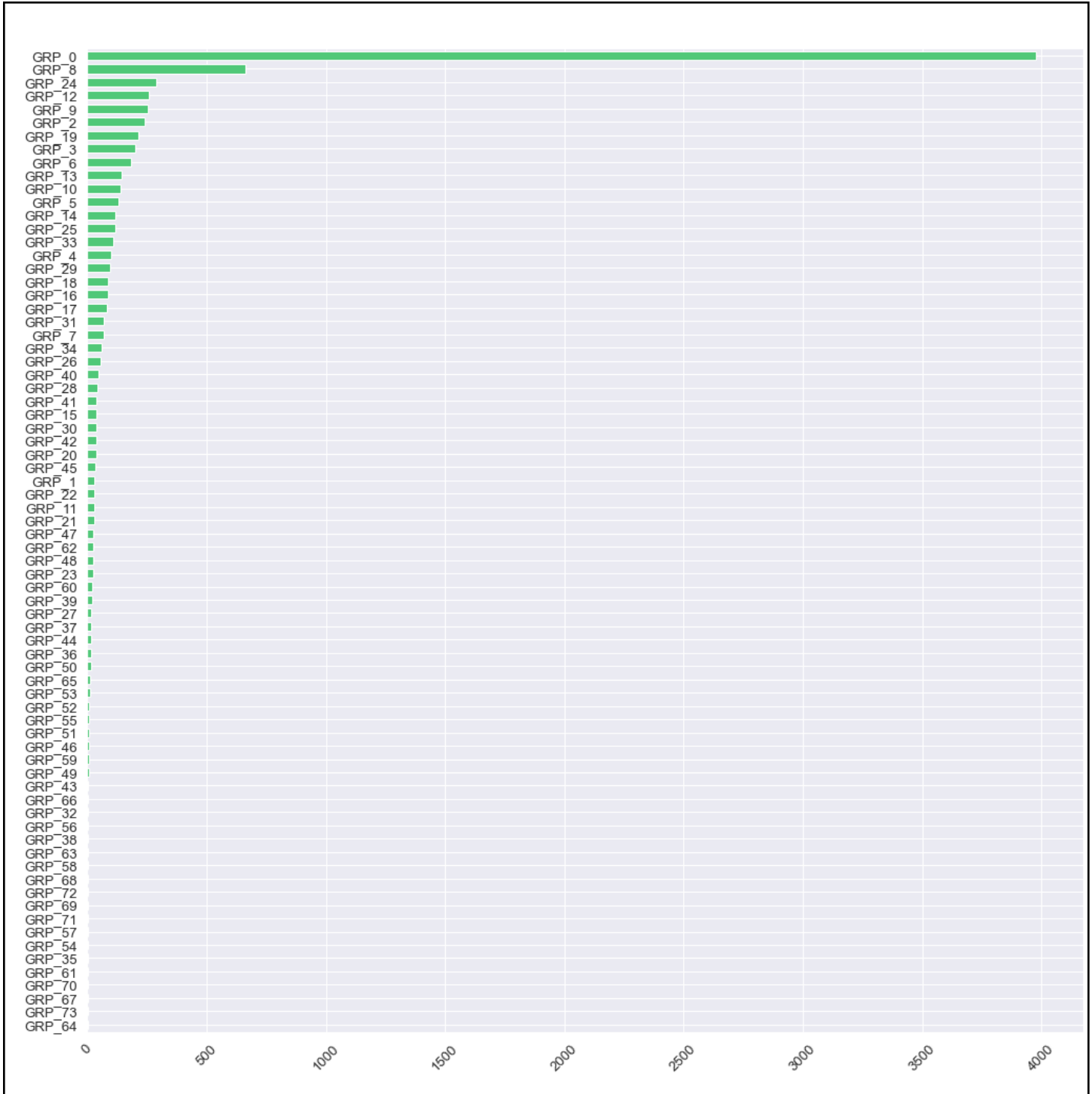
## 2. Summary of the approach to EDA and Pre-Processing

### 2.1. Project Life-Cycle



### 2.2.1. Target Distribution





- The Target class distribution is highly skewed and heavily imbalanced as most incidents are from Group 0, followed by Group 8, 24, 12, 9, 2. We find an imbalanced dataset for the rest of the groups.

- A large no. of entries for “Group 0”, which account for ~47% of the data and the remaining, are grouped as “Other” as there is not much information with the groups individually.

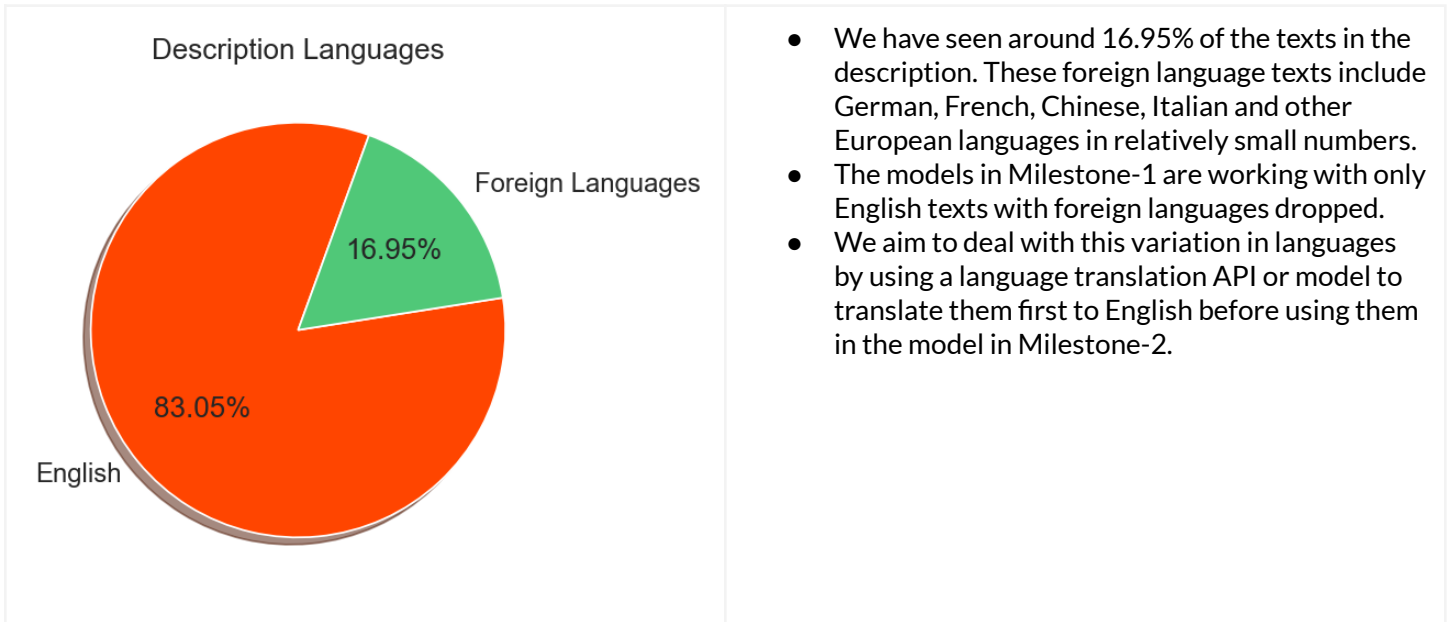
### 2.2.2. Choosing a Metric

- This is a multi-class classification problem, where the machine learning model will try to predict if each row is one of the 74 possibilities.
- The majority class is GRP\_0, which occurs in 46.78% of the observations.
- The most common metrics for a multi-class classification problem are AUC, F1-score and accuracy.
- Accuracy is not suitable for an imbalanced classification problem. (Note that a model that always predicts GRP\_0 will always get an accuracy of 46.78%)
- We will choose F1-Score if the majority class is more important than the smaller classes.
- We would choose AUC if we also care about the smaller classes.
- So, we use the AUC score to train models if there is an imbalance.
- Or, use accuracy if we fix the imbalance.

*As we want to classify the tickets into all functional groups and all functional groups are given equal importance, we choose **Accuracy** as the final metric to score model performance on the final balanced dataset.*

### 2.2.3. Language Detection





#### 2.2.4. Keyword Extraction

- [YAKE](#) is a lightweight unsupervised automatic keyword extraction using an unsupervised approach that rests on statistical text features extracted from single documents to select the most important keywords of a text and is independent of corpus, domain and language.
- The ten state-of-the-art unsupervised approaches followed here are TF.IDF, KP-Miner, RAKE, TextRank, SingleRank, ExpandRank, TopicRank, TopicalPageRank, PositionRank and MultipartiteRank.
- These keywords can then be used as a separate feature for the models and also for unsupervised clustering of groups based on the keywords later on.

```
print(test)
```

```
circuit outage: vogelfontein, south africa mpls circuit is down at 8:14 am et on 08/08
```

```
# !pip -q install yake
```

```
import yake
```

```
language = "en"
```

```
max_ngram_size = 5
```

```
duplication_threshold = 0.9
```

```
numOfKeywords = 1
```

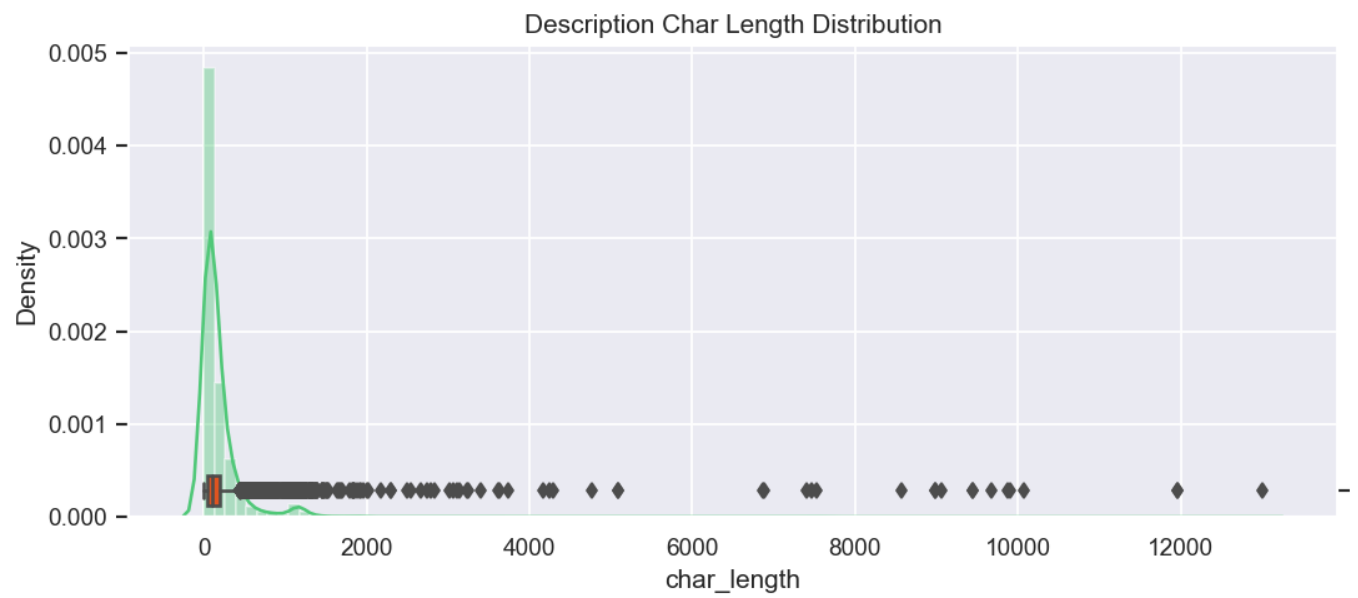
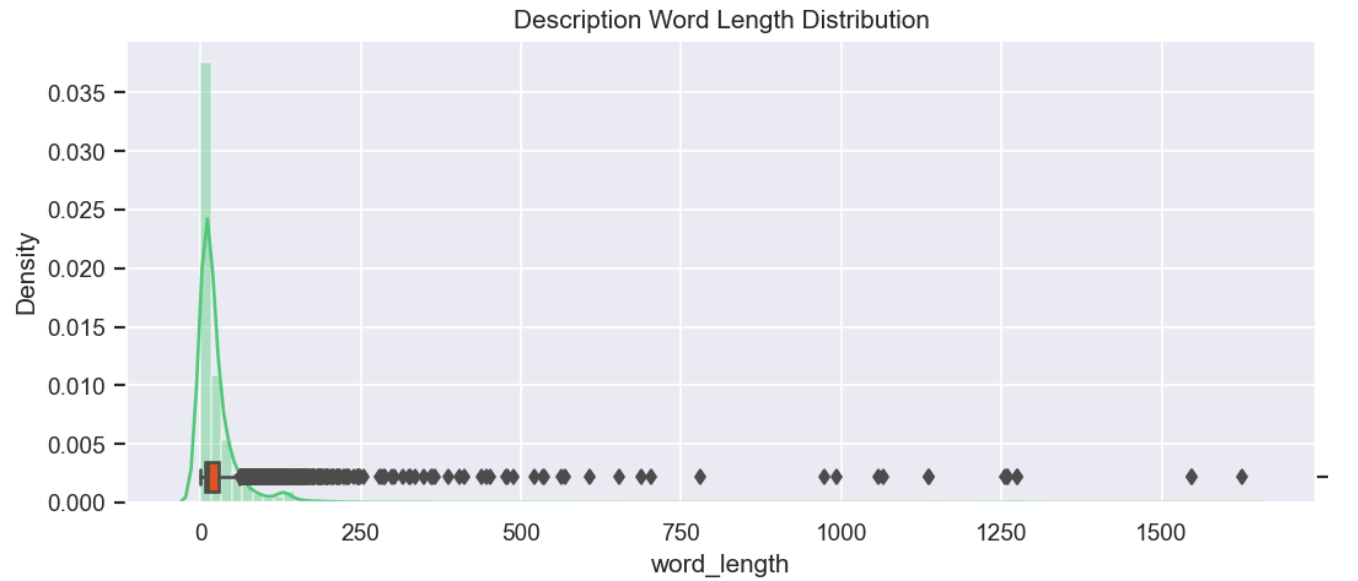
```
custom_kw_extractor = yake.KeywordExtractor(lan=language,  
                                             n=max_ngram_size,  
                                             dedupLim=duplication_threshold,  
                                             top=numOfKeywords,  
                                             features=None)
```

```
k = custom_kw_extractor.extract_keywords(test)
```

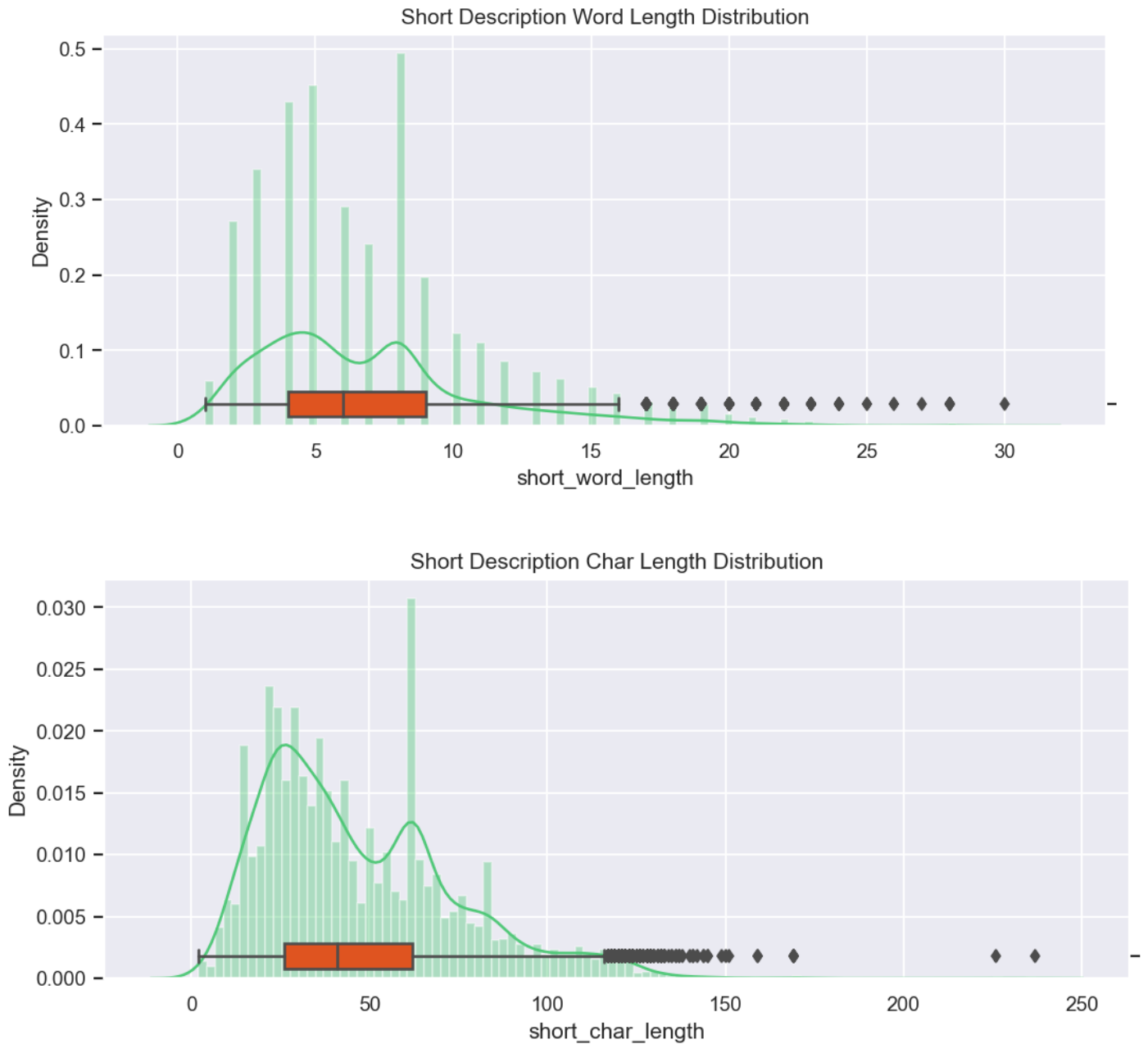
```
k[0][0]
```

```
'south africa mpls circuit'
```

## 2.2.5. Description Word and Character Counts distribution



## 2.2.6. Short Description Word and Character Counts Distribution



- Most descriptions are between 6 and 28 words long, with the median at 41 (106 characters) and the mean at 27.2 with relatively few outliers ranging to 1625 words.
- Most Short descriptions are between 4 and 9 words long, with the median at 6 (41 characters) and the mean at 6.92 with relatively few outliers ranging to 28 words.

## 2.2.7. Inconsistencies & Discrepancies found during EDA

Clean up the unwanted information from initial observations. Imputing the dataset which has no data, one and two-word length by their corresponding short description.

- No word length  $\Rightarrow$  Imputed the description with the corresponding short description.

short_description	description	caller	group	char_length	word_length	short_char_length	short_word_length
authorization add/delete members	\r\n\r\n	hpmwliog kqtnfvrl	GRP_0	5	0	33	3
browser issue :	\r\n	fgejnhux fnkymoh t	GRP_0	2	0	16	3

- One word length  $\Rightarrow$  Drop the row as the descriptions have no discernible information.

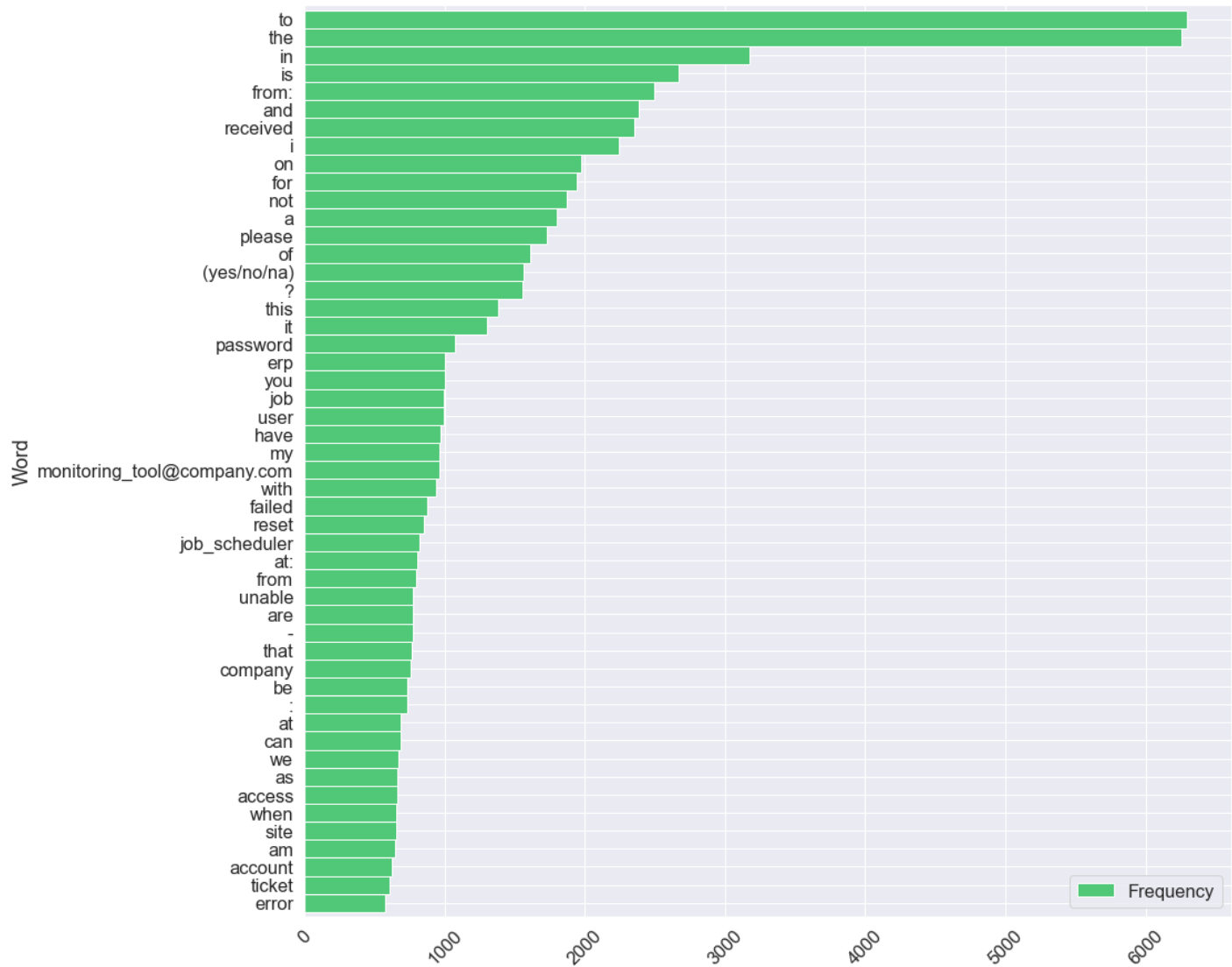
short_description	description	caller	group	char_length	word_length	short_char_length	short_word_length
s	s	gzjtweph mnsllwfqv	GRP_0	1	1	1	1

- Fix the encoding

short_description	description	caller	group	char_length	word_length	short_char_length	short_word_length
ç"µè,"ä, èf½å¼€æœ°	æ—©ä,Šä,Šç ç "µè,"æ%o"ä, å ¼€ä€,	mzerdt op xnlytcz j	GRP_30	30	1	18	1
outlookæ"¶å°°ç ®±ä, folderå ~ä, °æ- å¤©ä, €ä, ¢ ol...	outlookæ"¶å°°ç ç®±ä, folderå ~ ä, °æ- å¤©ä, €ä, ¢fol...	bxfdkio l mdqlsz vc	GRP_30	73	1	73	1

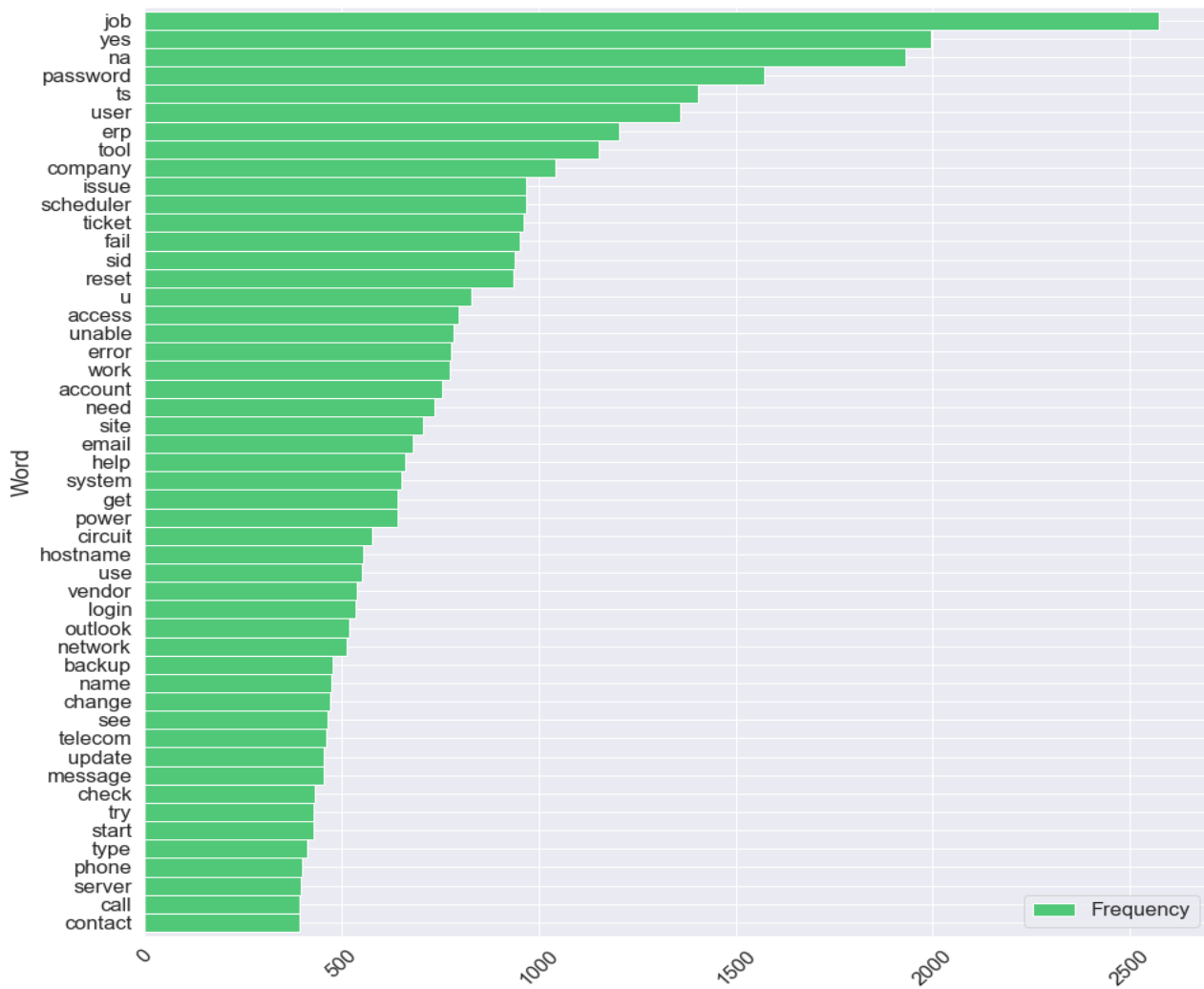
Description column and cleaned up data is generated for further analysis.

## 2.2.8 (a) Word Frequency Distributions



- We have observed that words like "to", "the", "in".. etc., are occurring most frequently in the descriptions. These words will not add any predictive power to the models and will need to be removed as part of the stopwords removal process during data pre-processing.
- Also, anchor words like "from:" and "received" and email addresses, punctuations, numbers are also occurring relatively frequently. We will remove these as well as part of the pre-processing.

## 2.2.8 (b) Word Frequency Distribution after Data Cleaning



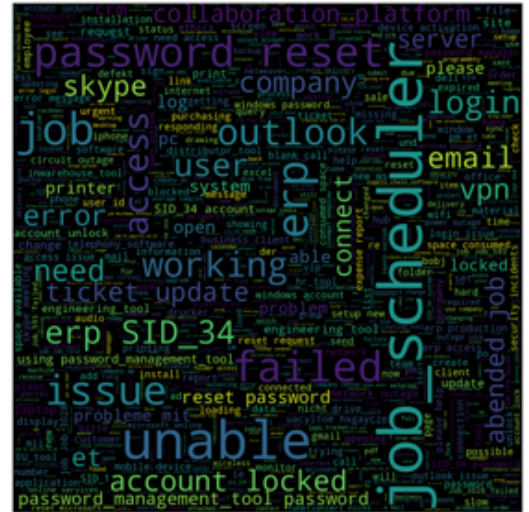
- The distribution of words is shown above after the data cleaning including the removal of stopwords, anchor words, numeric tokens, extra punctuation.
- Indicative words like "job", "password", "user" and "issue" are among the most frequent words.

### 2.2.9. Analysis using Word Clouds

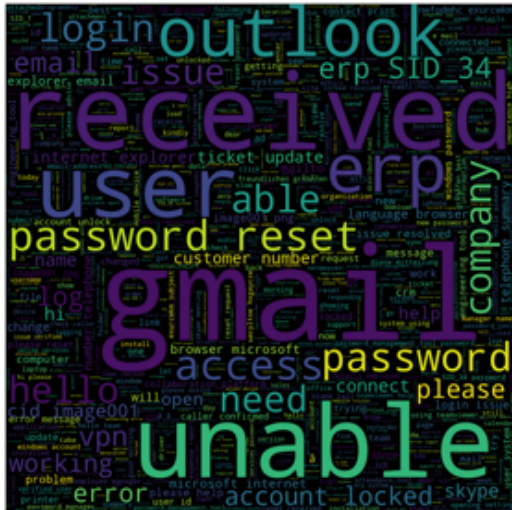
- Descriptions WordCloud



- Short Descriptions WordCloud



- Group 0 WordCloud



- Other Groups WordCloud



- WordCloud is a data visualization technique used for representing text data in which the size of each word indicates its frequency or relative importance within the dataset.
- Significant textual data points can be highlighted using a word cloud. WordClouds have been generated with All available words & top 50 words.
- We have also inferred a few observations over the target class – Assignment groups with word clouds for the top 50 words from each group.



## 2.2.10. Cleaning Outage Questionnaires

```
print(outage_df.description.tolist()[0])
```

```
what type of outage: ____network ____circuit ____x__power (please specify what type of outage)
1. top 23 cert site ? ____yes____ (yes/no/na)
2. when did it start ? ____4:31 pm et on 10/30. ____
3. scheduled maintenance ( power) ? __yno____ (yes/no/na) company power ____ provider power ____
4. scheduled maintenance ( network) ? __no____ (yes/no/na) company maint____ (yes/no) provider maint/ticket # ____
5. does site have a backup circuit ? __yes____ (yes/no/na)
6. backup circuit active ? ____na____ (yes/no/na)
7. site contact notified (phone/email) ? ____ (yes/no/na)
8. remote dial-in ? ____na____ (yes/no/na)
9. equipment reset ? ____na____ (yes/no/na)
10. verified site working on backup circuit ? __na____ (yes/no/na)
11. vendor ticket # ( global_telecom_1, verizon, telecom_vendor_1, telecom_vendor_2 ) ____global_telecom_1#000000223670658 ____
12. notified gsc ____ (yes/no/na) cert started ? ____ (yes/no/na)
13. additional diagnostics
```

```
print(clean_outageq(outage_df.description.tolist()[0]))
```

```
what type of outage:networkcircuitxpower please specify what type of outage top cert siteno when did it start pm et on scheduled maintenance
ket does site have a backup circuit backup circuit active nano site contact tified remote dialin na equipment reset na verified site working
om tifiedgsccert startedno additional diagstics
```

- These logs are related to outages in systems/products. These forms have a lot of noise in it.
- Here, we have extracted the required information which contains the response for the provided questions.
- Further, we clean up all the irrelevant information.

## 2.2.11. Cleaning Security Logs

```
# example string of a security log
test_string = '''
source ip: 10.16.90.249
source hostname: android-ba50a4497de455a
source port: 55198
source mac address: 50:2e:5c:f0:f6:98
system name :
user name:
location :
sep , sms status :
field sales user ( yes / no ) :
dsw event log:
-----
=====
event data
=====
related events:
event id: 84682727
event summary: internal outbreak for 137/udp
occurrence count: 505
event count: 1

host and connection information
source ip: 10.16.90.249
source hostname: android-ba50a4497de455a
source port: 55198
source mac address: 50:2e:5c:f0:f6:98
destination hostname: [no entry]
destination port: 137
connection directionality: internal
protocol: udp

device information
device ip: 80.71.06.702
device name: company-european-asa.company.com-1
log time: 2016-09-26 at 08:23:55 utc
action: blocked
cvss score: -1

scwx event processing information
sherlock rule id (sle): 537074
inspector rule id: 186739
inspector event id: 639601949
ontology id: 200020003203009162
event type id: 200020003203009062
agent id: 103761
event detail:
sep 26 08:23:55 80.71.06.702 %asa-4-106023: deny udp src inside:10.16.90.249/55198 dst noris:100.74.211.4/137 by access-group "acl_inside" [0x30e3d92a, 0x0]

[correlation_data]
sep 26 04:23:54 60.43.89.120 dhcpgd[23598]: dhcpack on 10.16.90.249 to 50:2e:5c:f0:f6:98 (android-ba50a4497de455a) via eth2 relay 10.16.88.2 lease-duration
691200 (renew)
sep 26 08:23:55 80.71.06.702 %asa-4-106023: deny udp src inside:10.16.90.249/55198 dst noris:100.74.211.1/137 by access-group "acl_inside" [0x30e3d92a, 0x0]

[correlation_data]
sep 26 04:23:54 60.43.89.120 dhcpgd[23598]: dhcpack on 10.16.90.249 to 50:2e:5c:f0:f6:98 (android-ba50a4497de455a) via eth2 relay 10.16.88.2 lease-duration
691200 (renew)
sep 26 08:23:55 80.71.06.702 %asa-4-106023: deny udp src inside:10.16.90.249/55198 dst noris:100.74.211.12/137 by access-group "acl_inside" [0x30e3d92a, 0x0]

[correlation_data]
sep 26 04:23:54 60.43.89.120 dhcpgd[23598]: dhcpack on 10.16.90.249 to 50:2e:5c:f0:f6:98 (android-ba50a4497de455a) via eth2 relay 10.16.88.2 lease-duration
691200 (renew)
sep 26 08:23:55 80.71.06.702 %asa-4-106023: deny udp src inside:10.16.90.249/55198 dst noris:100.74.211.14/137 by access-group "acl_inside" [0x30e3d92a, 0x0]

[correlation_data]
sep 26 04:23:54 60.43.89.120 dhcpgd[23598]: dhcpack on 10.16.90.249 to 50:2e:5c:f0:f6:98 (android-ba50a4497de455a) via eth2 relay 10.16.88.2 lease-duration
691200 (renew)

ascii packet(s):
[no entry]
hex packet(s):
[no entry]'''
```

```
>clean_sec_logs(test_string)

>after initial cleanup:
```

	cleaisource	androidbae	portsource	addressffsystem	name	user	name	location	sep	sms	status	field	sales	dsw	event	log	event	data	related
event idevent summary internal outbreak forudp occurrence countevent count host and connection information ipsourc androidbae portsource addressffdestination no entry destination portconnection directionality internal protocol udp device information device ipdevice name companyeuropeanasa.comany.comlog timeatut action blocked cvss score scwx event processing information sherlock rule id sleinspector rule idinspector event identityology idevent type idagent id event detail sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleasedurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleaseddurationrenew ascii packets no entry hex packets no entry																			

```
>after duplicates removal (final cleaned up string):

ipsource androidbae portsource addressffsystem name user location sep sms status field sales dsw event log data related idevent summary internal outbreak forudp occurrence countevent count host and connection information addressffdestination no entry destination portconnection directionality internal protocol udp device infoce companyeuropeanasa.comany.comlog timeatut action blocked cvss score scwx processing sherlock rule id sleinspector idnspector identityology type idagent detail sepa deny udp src insidedst norisby accessgroup acl_inside correlation_data sepdhcpd dhcpack ontoffandroidbae via ethrelayleasedurationrenew ascii packets hex"""
```

- In the security log clean-up, we have removed IP addresses, special characters, extra whitespace in the event data descriptions and duplicate entries.
- This function is specific to the security/event logs present in the dataset, which start with a specific pattern.

## 2.3 Key Insights and Takeaways

- 74 Assignment groups found - Target classes.
- **Group 0** is the majority class which accounts for ~47% of the data, and the remaining groups are relatively much less frequent, resulting in **highly imbalanced data**.
- Around 17% of the descriptions were found to be in **Non-English languages**.
- Several Emails were found in the description.
- Some descriptions have entire security/event logs.
- Symbols & other Non-ASCII characters were detected in the description.
- Hyperlinks, URLs, Email Addresses, Telephone Numbers & other irrelevant information was found in the descriptions.
- Blanks found either in the short description or description field.
- Few descriptions were the same as the short description.
- Few words were combined together.
- Spelling mistakes and typos were found in the data.
- Contraction words were found in the merged description and expanded for ease of word modelling.

## 2.4 Final Pre-Processing Techniques applied

```
print(test)
```

```
received from: tbvpkjoh.wnxzhqoa@gmail.com

i need access to the following path.  please see pmgzjiky potmrkxy for approval.

tbvpkjoh wnxzhqoa
company usa plant controller
tbvpkjoh.wnxzhqoa@gmail.com<tbvpkjoh.wnxzhqoa@gmail.com>

ticket update on inplant_872683
unable to login to collaboration_platform // password reset
all my calls to my ip phone are going to warehouse_toolmail, it is not even ringing.
sales area selection on opportunities not filtering to those in which the account
```

```
%%time
```

```
cleaned = preprocess_text(test)
pprint(cleaned, compact=True)
```

```
('need access follow path see pmgzjiky potmrkxy approval company usa plant '
'controller ticket update inplant 872683 unable login collaboration platform '
'password reset call ip phone go warehouse toolmail even ring sale area '
'selection opportunity filtering wch account')
```

```
Wall time: 49.5 ms
```

Below steps have been performed for initial pre-processing and clean-up of data in the preprocess\_text function:

- Fix text encoding using `ftfy.fix_text` — A lot of text in the data was being misinterpreted as some gibberish text (æ%“å¼€ outlook) when in reality they were Chinese characters (打开 outlook)
- Parse email messages to retain only subject and body — Parse the mails to strip out headers, salutations, attachments etc., to retain only the relevant message.
- Clean up emails, links, website links, telephone numbers — Strip out any of this unnecessary information using regex patterns.
- Clean up anchor words like: 'Received from:', 'name:', 'hello', 'hello team', 'cid'... etc., — Strip out any of these filler words which add no information to the model

- Clean up outage questionnaires — Cleans the outage questionnaires within the dataset by removing extra anchors question numbers answer options like (yes/no/na) etc., which add no relevant information.
- Clean up security logs — Clean the logs in the data by removing unnecessary information using regex patterns
- Clean HTML tags wherever they exist in the data
- Clean Blank (/r /n) characters
- Strip caller names in descriptions — Caller names were found to be present in the descriptions as well, these values were tokenized and stripped out if found in the descriptions.
- Translate/Normalize accented characters (á -> a)
- Convert Unicode characters to Ascii
- Expand contractions (they're -> they are)
- Clean stopwords & a few custom stopwords were found by analyzing the text
- Clean up extra whitespaces between words & Tokenize
- Remove gibberish — A lot of gibberish was still found to be in the text; this was stripped out using regex patterns
- Remove extra punctuation
- Changed the case sensitivity of words to lowercase
- Lemmatize the tokens in the final string
- Replaced Null values in Short description & description with space

### 3. Dealing with Class Imbalance Problem

#### 3.1. Clustering the Functional Groups based on frequency distribution

- Most common/simple issues are handled by L1 and L2 groups. More specialized issues are handled by L3.

```
def group_clustering(top_frequency_groups):
    descending_order = dataset['group'].value_counts().sort_values(ascending=False).head(top_frequency_groups)
    Cluster_1=[]
    Cluster_2=[]
    Cluster_1.extend(list(descending_order.index))
    Cluster_2.extend(list(set(dataset['group'].unique())-set(descending_order.index)))
    return Cluster_1,Cluster_2

# the value of top_frequency_groups=9 turns out to be best fit
L12,L3=group_clustering(top_frequency_groups=9)
print('L12')
pprint(L12, compact=True)

print('\nL3')
pprint(L3, compact=True)
```

```
L12
['GRP_0', 'GRP_8', 'GRP_24', 'GRP_12', 'GRP_9', 'GRP_2', 'GRP_19', 'GRP_3',
 'GRP_6']
```

```
L3
['GRP_53', 'GRP_35', 'GRP_37', 'GRP_58', 'GRP_73', 'GRP_25', 'GRP_23', 'GRP_42',
 'GRP_36', 'GRP_59', 'GRP_71', 'GRP_52', 'GRP_29', 'GRP_27', 'GRP_33', 'GRP_10',
 'GRP_20', 'GRP_43', 'GRP_31', 'GRP_44', 'GRP_55', 'GRP_68', 'GRP_15', 'GRP_63',
 'GRP_18', 'GRP_65', 'GRP_11', 'GRP_32', 'GRP_5', 'GRP_69', 'GRP_26', 'GRP_1',
 'GRP_62', 'GRP_70', 'GRP_46', 'GRP_13', 'GRP_39', 'GRP_49', 'GRP_22', 'GRP_60',
 'GRP_30', 'GRP_40', 'GRP_72', 'GRP_51', 'GRP_34', 'GRP_17', 'GRP_48', 'GRP_50',
 'GRP_57', 'GRP_38', 'GRP_45', 'GRP_61', 'GRP_16', 'GRP_41', 'GRP_54', 'GRP_14',
 'GRP_67', 'GRP_47', 'GRP_21', 'GRP_66', 'GRP_4', 'GRP_28', 'GRP_7', 'GRP_56',
 'GRP_64']
```

#### 3.2. Language Translation

The objective is to detect presence of Non-English languages in the dataset and using translation technique to preserve insights from non english text data

- **Detection of Non English Languages using fastText language identification model:**

- With the help of this pre-trained model, we were able to detect the presence of non english languages such as German, Mandarin, Portuguese, French etc.,

```
dataset.language.value_counts()
```

```
en      7962
de      484
zh       32
pt        8
fr        3
es        2
tl        2
fi        2
ca        2
pl        1
it        1
Name: language, dtype: int64
```

- We also measured the language confidence metrics to consider correctness of language detection with language confidence > 0.6, and the descriptions with language confidence < 0.6 are considered as english

```
# check some predictions with less confidence
dataset[(dataset.language != 'en') & (dataset.language_confidence > 0.6)][['merged_description', 'language',
```

	merged_description	language	language_confidence
223	probleme bluescreen hallo es ist erneut passie...	de	0.997537
251	reset password bitte passwort fr mail zurckset...	de	0.870327
255	probleme mit laufwerk laeusvjo	de	0.992616
265	hallo netweaver funktioniert nicht mehr bzw ka...	de	0.999658
270	neues passwort fur accountname tgryhu hgygrtui...	de	0.992012

- **Translating Non English Languages:**

- Translation using txtai pipelines which use hugging-face language translation models as backend
- The pipeline has logic to detect the input language, by loading the relevant model, it handles translating from source to # target language and return results.
- The translation pipeline also has built-in logic to handle splitting large text blocks into smaller sections the models can handle.



- **Sample Translations:**

```
# sample translations
translated_dataset[translated_dataset.language == 'de'].sample(20)[['merged_description', 'translated_description']]
```

	merged_description	translated_description
4397	block necessary unlock user ghjvreicj immediat...	block necessary unlock user ghjvreicj indirect...
1175	hallo zusammen bitte das iphone freischalten...	hello together please unlock the iphone for ma...

```
# sample translations
translated_dataset[translated_dataset.language == 'zh'].sample(20)[['merged_description', 'translated_description']]
```

	merged_description	translated_description
1953	笔记本重新装下系统 把我的笔记本重新装下系统	The notebooks reset the system, and my noteboo...
6534	网络不通 网络不通,右下角网络图标显示未连接到网络。	The network does not work. The network icon at...

### 3.3. Data Augmentation

- Data augmentation techniques are used to increase the amount of data by adding slightly modified copies of already existing data or newly created synthetic data from existing data. It acts as a regularizer and helps reduce overfitting when training a machine learning model
- We have implemented Data Augmentation with MarianMT using Back-Translation

#### Perform Augmentation using English <-> Spanish

```
en_texts = ['Cannot access website', 'I hated the food', "I can't login to my vpn"]
```

```
aug_texts = back_translate(en_texts, source_lang="en", target_lang="es", verbose=True)
print(aug_texts)
```

Intermediate Target Language texts:

```
['No se puede acceder al sitio web', 'Odiaba la comida.', 'No puedo acceder a mi vpn']
['Cannot access the website', 'I hated food.', "I can't access my vpn"]
```

### 3.4. Upsampling

- Finally, we upsample the minority class to balance the target distribution so that model learns better on both classes.

```
augmented_df.label.value_counts()
```

```
0    13109  
1     9716  
Name: label, dtype: int64
```

```
upsampled_df.label.value_counts() # BALANCED DATA!
```

```
1    13109  
0    13109  
Name: label, dtype: int64
```

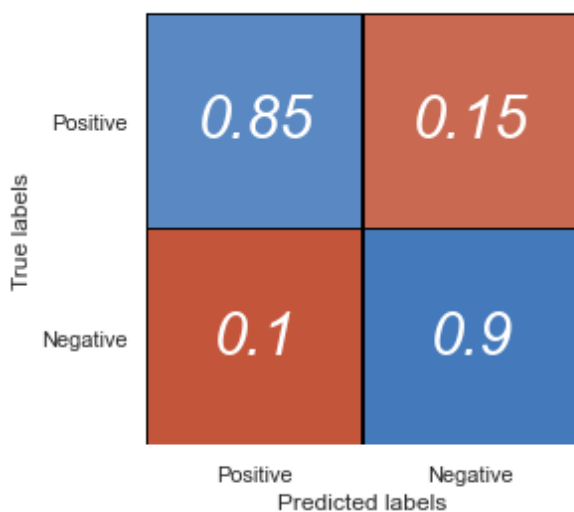
## 4. Model Building

### 4.1. Machine Learning Models

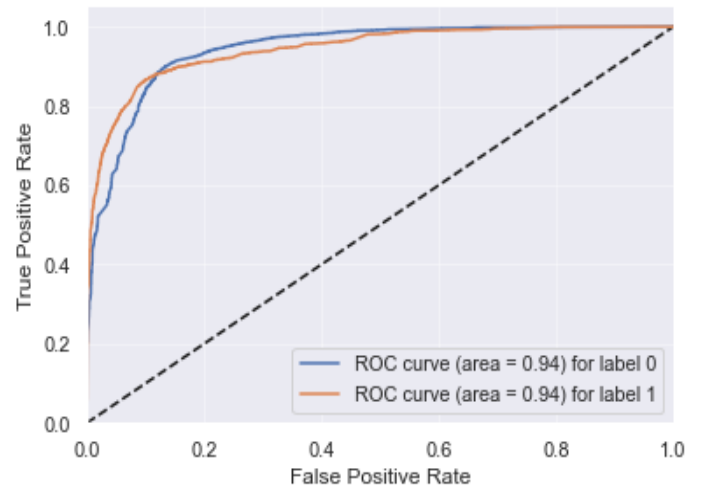
- The machine learning models were built based on random search cross-validation, two separate data sets were used for model building.
- The training dataset for the ML models consisted of 2 target classes: **L12 and L3**.
- Metrics for the models:

Model	Test Accuracy
LightGBM Classifier	71.41%
XGBoost Classifier	83.97%
Random Forest Classifier	<b>87.19%</b>

Confusion matrix (Random Forest Classifier)



ROC Curve (Random Forest Classifier)



## 4.2. Deep Learning Models

Model	Test Accuracy
<b>All Classes Classification (Functional groups)</b>	
<b>Bi-Directional LSTM</b>	65.87%
<b>Tfidf Vectors + Feature Selection + Feed-forward Neural Net</b>	66.80%
<b>Binary Classifiers (L12 - L3 groups)</b>	
<b>Convolution Blocks (Dimensionality Reduction) + Bi-LSTM</b>	82.67%
<b>Stratified KFold Validation + Tfidf Vectors + Feature Selection + ANN</b>	83.60%
<b>Foreign Langs. Translated + Generated Keywords Feature + Stratified KFold Validation + Tfidf Vectors + Feature Selection + ANN</b>	91.21%
<b>Foreign Langs. Translated + Data Augmentation + Stratified KFold Validation + Tfidf Vectors + Feature Selection + ANN</b>	91.40%
<b>Foreign Langs. Translated + Data Augmentation + Upsample (Balanced) + Stratified KFold Validation + Tfidf Vectors + Feature Selection + ANN</b>	<b>92.56%</b>

- We have observed that the Bidirectional LSTM model is performing much better relative to the LSTM model on most of the samples. One reason behind this could be that the bidirectional model takes into account the past and the future context of a sequence and is hence more robust in dealing with the noise that might be biasing the vanilla LSTM model and also understanding the context of the words in a description.
- The best model was trained using the preprocessed dataset where we've translated non-english texts, used back translation to do data augmentation for both the minority (L12) and majority (L3) groups to help deep learning models learn the significant keywords better.
- This was followed by upsampling of the minority class after data augmentation to balance the target distribution.
- Finally, The custom model with TF-IDF vectorization and feature selection approach worked best for this project. Here, we are processing more contextual information as compared to the generalized nature of word embeddings and through Feature Selection, we are able to capture the most important features and avoid noise going into the model.

## 4. Key Learnings & Further Improvements

### 4.1. Learnings:

- The dataset is highly imbalanced, which could be an inherent limitation that affects the performance of our classification model.
- Oversampling of minority classes with text augmentation was used to handle the imbalance in the data after data augmentation.
- The presence of other languages in the dataset is an inherent limitation within the dataset, which will limit the language models to learn properly. This is further complicated by the fact that we may receive similar foreign language descriptions in out-of-sample data which have to be translated first adding to the latency of the inference pipeline.
- Language Translation was done using Hugging Face pre-trained translation models.
- Even though we found that the caller column had a significant correlation with the target column by performing a  $\chi^2$  test for two distributions, we have decided to drop the caller column as an input feature to the models as considered in the tradeoff:
  - New queries from an old caller could give some prior probabilities/info to the model.
  - But it wouldn't help on new callers in out-of-sample data as the caller doesn't really indicate what their ticket or issue is. So, Real-world performance will degrade.
  - In view of the above tradeoff, we have dropped the caller id column as a feature.
- As Deep Learning requires a large set of data to work well, the limited number of data points within our dataset is an inherent limitation to deep learning models. To be able to train production-ready models for our use-case. We would need to gather more data/issues on some of the minority classes.

## 4.2. Further Improvements:

- Gather more data from external sources for the issues so that our training dataset target distribution is better.
- Stacked model on top of the binary classifier which further classifies among the other groups.
- Transfer Learning using some pre-trained NLP models to utilize the language understanding of these models and fine-tune them for ticket classification.
- Add Attention Layer to the model architecture to capture context and significant keywords in the descriptions better.
- Use other pre-trained embeddings like Fasttext, BERT, ELMO, Flair Word Embeddings to capture the context/meaning of the natural language words.
- Train transformer based model architectures (e.g, BERT).
- Productionise the model by deploying locally or on cloud.
- An active learning approach might also be implemented, which can learn as new types of issues come up through to customer service.