

AIML Online Capstone Project
AUTOMATIC TICKET ASSIGNMENT
INTERIM REPORT

Table of Contents

Real Problem	3
Objective.....	3
Summary of the problem statement, data, and findings:	3
Observations of the dataset:	4
Approach to Data Cleansing, EDA and Text Pre-processing:	4
a) Data Cleaning.....	4
b) EDA.....	7
c) Text-Pre-processing.....	9
Lemmatization & Stop words removal	9
Analysis using WordCloud	10
Model Selection.....	11
Model Performance	11
How to improve your model performance?	14

Real Problem

In any IT industry, Incident Management plays an important role in delivering a quality support to customers. An incident ticket is created by various groups of people within the organization to resolve an issue as quickly as possible based on its severity. Whenever an incident is created, it reaches the Service desk team, and then it gets assigned to the respective teams to work on the incident. The Service Desk team (L1/L2) will perform basic analysis on the user's requirement, identify the issue based on given descriptions and assign it to the respective teams.

The manual assignment of these incidents might have the below disadvantages:

- More resource usage and expenses.
- Human errors - Incidents get assigned to the wrong assignment groups
- Delay in assigning the tickets
- More resolution times
- If a particular ticket takes more time in analysis, other productive tasks get affected for the Service Desk

Objective

From the given problem description, we could see that the existing system is able to assign 75% of the tickets correctly. So, our objective here is to build an AI-based classifier model to assign the tickets to right functional groups by analyzing the given description with an accuracy of at least 80%.

Summary of the problem statement, data, and findings:

We have received unstructured data as our dataset. The dataset has 8500 rows and 4 columns.

The 4 columns are:

1. **Short Description:** This column gives us a glance at the broad categorization of the complaint. Example: login issue, outlook, can't log in to VPN, unable to access hr_tool page, skype error
2. **Description:** This column gives us a little more detailed insight into the complaint that the end-user has. It talks about a brief description of the complaint/issue the end-user is facing. For example, my meetings/skype meetings are not appearing in my outlook calendar; can somebody please advise how to correct this?
3. **Caller:** This column has the end user's name who is filing the complaint or raising the request. Example: spxjnwir pjlcoqds, hmjdrvpb komuaywn
4. **Assignment Group:** This column talks about the category/group into which the complaint is classified for it to be directed to the right department for the issue to be resolved. Example: GRP_0, GRP_1

Observations of the dataset:

- The total number of incidents reported in the dataset is 8500.
- Caller names in a random fashion (may not be useful for training data)
- The dataset has English and German Language words.
- Email/chat format found in the description.
- There are a total number of 8 null records in the Short Description column and there is 1 null record in the Description column. There are no null records in the Caller and Assignment Group columns.
- There are a total of 74 different groups for the Assignment Group column.
- Approximately 50% of the dataset comprises complaints that are corresponding to GRP_0.
- There are a few rows of data that have the same text for the Short Description and the Description column.
- Few words were combined.
- Spelling mistakes and typo errors are found.

Approach to Data Cleansing, EDA and Text Pre-processing:

a) Data Cleaning

We performed the stated actions

- We started exploratory data cleaning by getting the basic information of our dataset.

```
<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
```

- Dropping the Callers Column from Dataset

```
# Removing the Callers column
callers = raw_data['Caller'].unique()
raw_data.drop(columns='Caller', inplace=True)
raw_data.head()
```

	Short description	Description	Assignment group
0	login issue	-verified user details.(employee# & manager na...	GRP_0
1	outlook	\r\n\r\nreceived from: hmjdrvpb.komuaywn@gmail...	GRP_0
2	cant log in to vpn	\r\n\r\nreceived from: eylqgodm.ybqkwiam@gmail...	GRP_0
3	unable to access hr_tool page	unable to access hr_tool page	GRP_0
4	skype error	skype error	GRP_0

- Removal of special characters, trailing spaces, numbers
- Converting strings to lower case

```
def clean_data(text):
    text = text.lower()
    text = ' '.join([w for w in text.split() if not is_valid_date(w)])
    text = re.sub(r"received from:", ' ', text)
    text = re.sub(r"from:", ' ', text)
    text = re.sub(r"to:", ' ', text)
    text = re.sub(r"subject:", ' ', text)
    text = re.sub(r"sent:", ' ', text)
    text = re.sub(r"ic:", ' ', text)
    text = re.sub(r"cc:", ' ', text)
    text = re.sub(r"bcc:", ' ', text)
    #Remove email
    text = re.sub(r'\S*@\S*\s?', '', text)
    # Remove numbers
    text = re.sub(r'\d+', '', text)
    # Remove new line characters
    text = re.sub(r'\n', ' ', text)
    # Remove hashtag while keeping hashtag text
    text = re.sub(r'#', '', text)
    text = re.sub(r'&?', 'and', text)
    # Remove HTML special entities (e.g. & amp;)
    text = re.sub(r'&\w*;', '', text)
    # Remove hyperlinks
    text = re.sub(r'https?:\V\/*\V\w*', '', text)
    # Remove characters beyond Readable format by Unicode:
    text = ''.join(c for c in text if c <= '\uFFFF')
    text = text.strip()
    # Remove unreadable characters (also extra spaces)
    text = ' '.join(re.sub("[^\u0030-\u0039\u0041-\u005a\u0061-\u007a]", " ", text).split())
    for name in callers:
        namelist = [part for part in name.split()]
        for namepart in namelist:
            text = text.replace(namepart, '')

    text = re.sub(r"\s+[a-zA-Z]\s+", ' ', text)
    text = re.sub(' +', ' ', text)
    text = text.strip()
    return text
```

```
# Apply the cleaning function to entire dataset
raw_data['Description'] = raw_data['Description'].apply(clean_data)
```

- Merging small groups

```
# Group the tickets count < 100 as GRP_A
Ticket1 = pd.DataFrame(raw_data['Assignment group'].value_counts())
Ticket1 = Ticket1.T
Ticket1
```

- Merging both Description Columns

```
# Merging "Short description" and "Description" column
raw_data['Description'] = raw_data['Short description'] + ' ' + raw_data['Description']
raw_data.drop(columns=['Short description', 'Count'], inplace=True)
raw_data.head()
```

	Description	Assignment group
0	login issue verified user details employee and...	GRP_0
1	outlook hello team my meetings skype meetings ...	GRP_0
2	cant log in to vpn hi cannot log on to vpn best	GRP_0
3	unable to access hr_tool page unable to access...	GRP_0
4	skype error skype error	GRP_0

Observations:

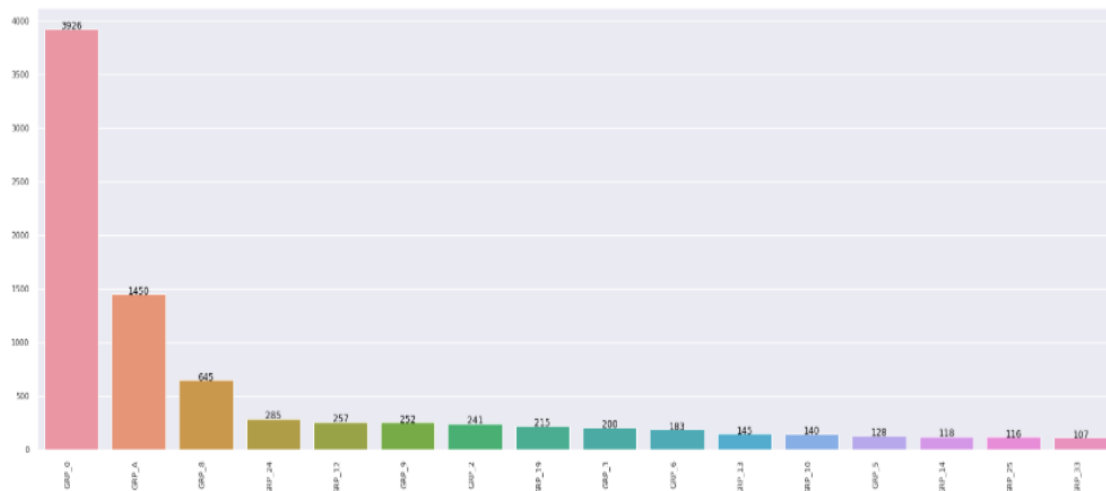
- Entire dataset is converted into lower case.
- Users email addresses will add NO value to our analysis, despite the fact that user id is given in the caller column. So, all email addresses are removed from the dataset
- All numerals are removed because they were dominating the dataset if we were converting them into their word representation otherwise.
- All punctuation marks are removed which used to be a hindrance in lemmatization.
- All occurrences of more than one blank spaces, horizontal tab spaces, new line breaks etc. have been replaced with single blank space.

b) EDA

- We merged all the small groups into one group. We also converted the entire dataset into lower case. Here is the new information on the new dataset.

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 8408 entries, 0 to 8499
Data columns (total 2 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Description      8408 non-null   object
1   Assignment group  8408 non-null   object
dtypes: object(2)
memory usage: 197.1+ KB
```

- The assignment group was analyzed again and it was 16 groups.
- As we observe in the below bar plot, the new dataset is extremely right-skewed. GRP_0, GRP_A, and GRP_8 together make to around 70% of the available data.



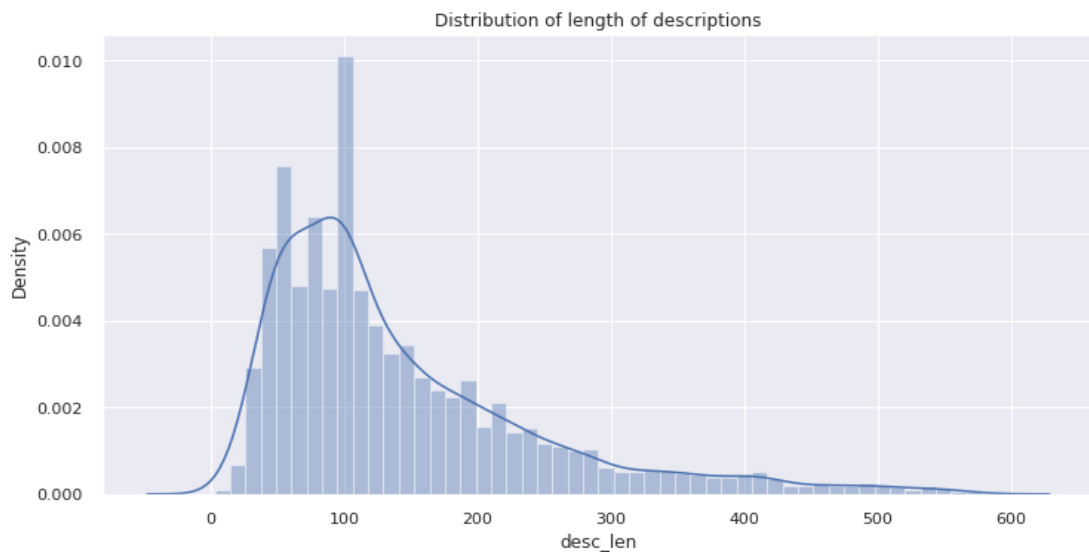
Observation:

- The Target class distribution is extremely right skewed.
- GRP_0 has max occurrences (amounting to 3926) which makes for almost ~50% of the data.
- The merged entries GRP_A has around 1450 occurrences.
- GRP_0, GRP_A and GRP_8 together makes to around 70% of the available data.

```
raw_data['desc_len'].describe(percentiles=[0.25,0.5,0.75,0.9,0.95])
```

```
count      8417.000000
mean       190.084828
std        319.030581
min         3.000000
25%        74.000000
50%       114.000000
75%       202.000000
90%       360.000000
95%       580.000000
max       6300.000000
Name: desc_len, dtype: float64
```

- We then figured out the length of each description. Here is the analysis.



Summary of the EDA:

- There is no clear pattern visible to describe the relation between group and description lengths.
- It is however visible that groups have outliers and some groups fall in the lower range while most of them have lengths in the range of 100 to 400.
- Half of the records have less than 120 words.
- 90% of descriptions have less than 620 words.
- The remaining 10% records are very long and extremely skewed in comparison with the remaining portion.

c) Text-Pre-processing

Lemmatization & Stop words removal

Stop words have been removed using nltk corpus modules.

Lemmatization is the process of grouping together the different inflected forms of a word so they can be analysed as a single item. Lemmatization is similar to Stemming but it brings context to the words. So, it links words with similar meanings to one word.

Here we have preferred Lemmatization over Stemming because lemmatization does morphological analysis of the words.

```
from nltk.corpus import stopwords
import nltk
nltk.download('wordnet')
nltk.download('punkt')
nltk.download('averaged_perceptron_tagger')

import spacy
nlp = spacy.load('en', disable=['parser', 'ner'])
allowed_postags=['NOUN', 'ADJ', 'VERB', 'ADV']
def lemmatize_text(text):
    doc = nlp(text)
    return ' '.join([token.lemma_ for token in doc if token.lemma_ != '-PRON-'])

raw_data['Description'] = raw_data['Description'].apply(lemmatize_text)
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Unzipping corpora/wordnet.zip.
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger.zip.
```

- We removed all the stop words so as to better understand our description and n-gram analysis.

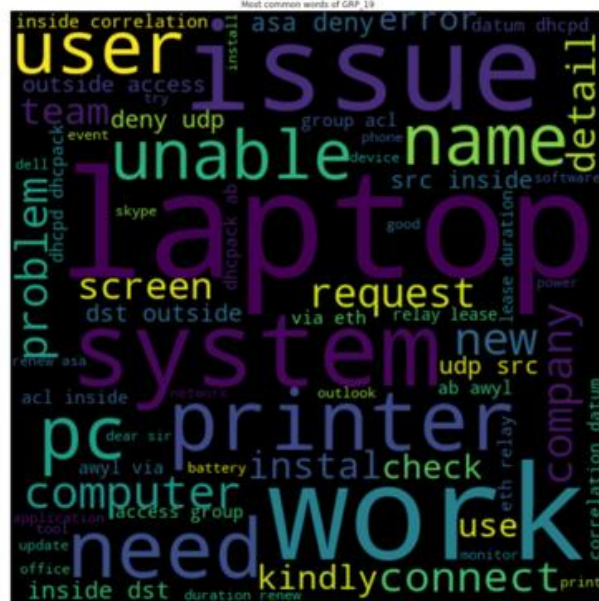
```
from wordcloud import STOPWORDS
from sklearn.feature_extraction.text import CountVectorizer
# Extend the English Stop Words
STOP_WORDS = STOPWORDS.union({'yes', 'na', 'hi', 'etc',
                              'receive', 'hello',
                              'regards', 'thanks',
                              'from', 'greeting',
                              'forward', 'reply',
                              'will', 'please',
                              'see', 'help', 'able'})

# Generic function to derive top N n-grams from the corpus
def get_top_n_ngrams(corpus, top_n=None, ngram_range=(1,1), stopwords=None):
    vec = CountVectorizer(ngram_range=ngram_range,
                        stop_words=stopwords).fit(corpus)
    bag_of_words = vec.transform(corpus)
    sum_words = bag_of_words.sum(axis=0)
    words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
    words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
    return words_freq[:top_n]
```

Analysis using WordCloud

Word cloud is a collection, or cluster, of words depicted in different sizes. The bigger and bolder the word appears, the more often it's mentioned within a given text and the more important it is.

Also known as tag clouds or text clouds, these are ideal ways to pull out the most pertinent parts of textual data, often also help business users compare and contrast two different pieces of text to find the wording similarities between the two.



Summary of the Text-Preprocessing:

It's indicative from the n-gram analysis and the word cloud is that the entire dataset speaks more about issues around

- password reset
- fail job & scheduler

Analysis on GRP_0 which is the most frequent group to assign a ticket to reveals that this group deals with mostly the maintenance problems such as password reset, account lock, login issue, ticket update etc.

Model Selection

- Built Unigram, Bi-gram and Tri-gram models from the bag of words
- Multinomial Naive Bayes
- K Nearest Neighbor (KNN)
- Support Vector Machine
- Decision Tree
- Random Forest
- Bidirectional LSTM

Model Performance

We tried to implement all these models to the data but the result we got very inaccurate. All the models are suffering from low test accuracy.

Naive Bayes

```
run_classification(MultinomialNB(), X_train, X_test, y_train, y_test)
```

```
Prediction Model: MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)
```

```
-----  
Training accuracy: 61.18%
```

```
Testing accuracy: 57.36%  
-----  
-----
```

K-Nearest Neighbor (KNN)

```
run_classification(KNeighborsClassifier(), X_train, X_test, y_train, y_test)
```

```
Prediction Model: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                                       metric_params=None, n_jobs=None, n_neighbors=5, p=2,
                                       weights='uniform')
```

```
-----  
Training accuracy: 76.33%
```

```
Testing accuracy: 68.76%  
-----  
-----
```

Support Vector Machine (SVM)

```
# SVM with Linear kernel
```

```
run_classification(LinearSVC(), X_train, X_test, y_train, y_test)
```

```
Prediction Model: LinearSVC(C=1.0, class_weight=None, dual=True, fit_intercept=True,
                             intercept_scaling=1, loss='squared_hinge', max_iter=1000,
                             multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
                             verbose=0)
```

```
-----  
Training accuracy: 95.78%
```

```
Testing accuracy: 73.22%  
-----  
-----
```

Decision Trees

```
run_classification(DecisionTreeClassifier(), X_train, X_test, y_train, y_test)
```

```
Prediction Model: DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                                          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')
```

```
-----  
Training accuracy: 99.69%
```

```
Testing accuracy: 61.52%  
-----  
-----
```

Random Forest Classifier

```
from sklearn.ensemble import RandomForestClassifier
run_classification(RandomForestClassifier(n_estimators=100), X_train, X_test, y_train, y_test)
```

```
Prediction Model: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                                         criterion='gini', max_depth=None, max_features='auto',
                                         max_leaf_nodes=None, max_samples=None,
                                         min_impurity_decrease=0.0, min_impurity_split=None,
                                         min_samples_leaf=1, min_samples_split=2,
                                         min_weight_fraction_leaf=0.0, n_estimators=100,
                                         n_jobs=None, oob_score=False, random_state=None,
                                         verbose=0, warm_start=False)
```

Training accuracy: 99.67%

Testing accuracy: 66.92%

Bidirectional LSTM Classifier

```
run_classification(model, train_padded, validation_padded, training_label_seq, validation_label_seq, arch_name=None, pipelineRequired=False, isDeepModel=True)
```

```
Epoch 1/10
53/53 - 29s - loss: 2.0492 - accuracy: 0.4759 - val_loss: 1.6388 - val_accuracy: 0.5493
Epoch 2/10
53/53 - 24s - loss: 1.4142 - accuracy: 0.6040 - val_loss: 1.4163 - val_accuracy: 0.5612
Epoch 3/10
53/53 - 24s - loss: 1.2355 - accuracy: 0.6446 - val_loss: 1.3826 - val_accuracy: 0.5689
Epoch 4/10
53/53 - 24s - loss: 1.1191 - accuracy: 0.6768 - val_loss: 1.3689 - val_accuracy: 0.5873
Epoch 5/10
53/53 - 25s - loss: 1.0010 - accuracy: 0.6992 - val_loss: 1.3608 - val_accuracy: 0.5730
Epoch 6/10
53/53 - 25s - loss: 0.8877 - accuracy: 0.7297 - val_loss: 1.3951 - val_accuracy: 0.5808
Epoch 7/10
53/53 - 25s - loss: 0.7757 - accuracy: 0.7598 - val_loss: 1.3764 - val_accuracy: 0.6010
Epoch 8/10
53/53 - 25s - loss: 0.6611 - accuracy: 0.8031 - val_loss: 1.4516 - val_accuracy: 0.5992
Epoch 9/10
53/53 - 25s - loss: 0.5831 - accuracy: 0.8264 - val_loss: 1.5169 - val_accuracy: 0.6134
Epoch 10/10
53/53 - 25s - loss: 0.5019 - accuracy: 0.8542 - val_loss: 1.5062 - val_accuracy: 0.6217
Prediction Model: <tensorflow.python.keras.engine.sequential.Sequential object at 0x7f2b53bda810>
-----
Training accuracy: 87.64%
Testing accuracy: 62.17%
-----
```

The summary of the results we got are attached below:

Model	Training Accuracy	Test Accuracy
Multinomial Naïve Bayes	61%	59%
K-Nearest Neighbour	76%	69%
Linear SVM	96%	73%
Decision Tree	99%	60%
Random Forest	99%	68%
Bidirectional LSTM	91%	60%

How to improve your model performance?

In all the models we have tried, the accuracy of each of the Statistical models is overfitted to a higher degree. One obvious reason is the dataset is highly imbalanced.

For milestone 2 we will try to train and tune the data in a more accurate way and we will deal with the imbalanced data set.