# Internship Report on Machine Learning and Deep Learning
# (Natural Language Processing)

Submitted by

**Udit Nath**
B.Tech in Computer Science and Engineering
Central Institute of Technology Kokrajhar

Under the guidance of

**Dr. Pradeep Kumar Roy**
Computer Science and Engineering Assistant Professor (Grade-II,Level-10)
IIIT Surat

Department of Computer Science and Engineering
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY, SURAT

**Abstract**

This is an internship report on Machine Learning and Deep learning introduction. This report contains the details of the various tasks performed during two months of internship and is compiled for recognition for the completion of the internship. Tasks such as data cleaning, data processing and featured engineering, Feature Engineering using Countvectorizer and Bag of Words approach, Feature Engineering using Tf-Idf, Applying Machine learning algorithms on the processed dataset, Implementation of CNN algorithm on the same dataset, Text Normalization and finally Applying ML and CNN models on Bi-lingual dataset and comparing their outcomes.

# Contents

# Chapter 1

# Introduction

This report contains the detailed analysis of the two months of internship in Machine Learning and Deep Learning in IIIT Surat.

It starts with the introductory data cleaning process. Omitting Null spaces and deleting stopwords are the basics which are to be mastered first. Later on the data cleaning process involves the noise reduction in the data. It includes, removing HTML tags, punctuation and unnecessary obstructions to data reading process.

Introduction of Bag-of-words, Count vectorizer, TF-IDF had been covered in the next part. Convolutional Neural Network is next covered for the 1D layers for text analysis. At the end, a simple Codemix Dataset of Bi-Lingual data is been process with the techniques that were used for the prior tasks.

# Chapter 2

# Work Done

This Chapter introduces the various methods used in the Machine Learning and Deep Learning processes and the detailed analysis of the same.

## 2.1 Experimental setup

### 2.1.1 Softwares used

Google Collab is used as a platform for the internship tasks and some words had been done on Jupyter NB which had been later shited in the Google Collab notebook.

| Python Libraries | Usage |
|---|---|
| Pandas | Used for data analysis and machine learning tasks |
| NumPy | Provides support for multi-dimensional arrays |
| Re | A regular expression (RE) is a language for specifying text search strings |
| NLTK | Making natural human language usable by computer programs |
| keras | Creating deep models |
| sklearn | Provides tools for machine learning and statistical modeling |

### 2.1.2 Hardware used

A portable system can be used. In this particular case Macbook Pro has been used containing Apple M1 chipset and 8-core CPU with 6 performance cores and 2 efficiency cores, 14-core GPU, 16-core Neural Engine, 200GB/s memory bandwidth.

## 2.2 Detailed Description

### 2.2.1 Task 1

The Task 1 contains the simple data cleaning process. For this task, a NPL disaster dataset has been used.

Firstly, the data is loaded with the help of Pandas library. Next, we check the null values present in the dataset. For this process, many iterations had been tried on the same dataset to experiment with the difference in results. Then the null value were filled with custom values which were user defined. inplace argument in the .dropna() function had been learnt as a bonus point for the further use. The use of frequent stopwords is repetitive in the dataset. So, it has to be removed in the cleaning process.



Figure 2.1: Text Cleaning processes

### 2.2.2 Task 2

For performing Task 2, the same dataset has been used. In this task the basics of feature engineering are known. Some preprocessing techniques like

cleaning the stop words had been implemented in this task. First of all we perform the basic data preprocessing tasks on the data. Then proceeded to data vsualization. We plot different graphs to see different outcome of the same dataset.Next we handle the co-relation in the dataset. At the end we detect the outliner using Interquartile range (IQR) and remove them.



Figure 2.2: Preprocessing of the data

### 2.2.3   Task 3

In this part of feature extraction, the report deals with various methods and techniques that are used for feature extraction. The theories and usage of Bag-of-words, Count Vectorizer, TF-IDF had been learnet during the task 3.

A bag-of-words is a representation of text that describes the occurrence of words within a document. It involves two things:
  A vocabulary of the known words.
  A measure of the presence of known words.

Example:
This is my pen.
This is where I go to study.

Bag of Words for sentence 1:
"This" = 1
"is" = 1
"my" = 1
"pen" = 1
"where" = 0
"I" = 0
"go" = 0
"to" = 0
"school" = 0

Limitations:
Vocabulary, Sparsity and Meaning of sentences

CountVectorizer is a method to convert text to numerical data.

CountVectorizer(parameters):

- lowercase = True (default)

- stop_words: Words that are insignificant and can be avoided

- stop_words = "english" — sklearn built in stop words list

- max_df=0.75 some words crosses the threshold of 75% and removed from the sparse matrix

- min_df=0.50 percentage of documents(0.50, ignore words appearing in 50% of documents)

- max_features = 3 it will select the 3 most common words in the data. binary = True the CountVectorizer no more takes into consideration the frequency of the term/word

- .fit_transform(df) = transforms the data to vector form

- .get_feature_names() = Learn the vocabulary dictionary and return document-term matrix

- .vocabulary = represents the position of the words in the matrix

### 2.2.4   Task 4

TF-IDF: Rescale the frequency of words by how often they appear in all documents, so that the scores for frequent words like "the" that are also frequent across all documents are penalized.

- Term Frequency: is a scoring of the frequency of the word in the current document

- Inverse Document Frequency: is a scoring of how rare the word is across documents



Figure 2.3: CountVectorizer and TF-IDF

### 2.2.5   Task 5

In this section, Machine Learning algorithms had been used for the same input of the previous tasks to compare the output of different approach.
Machine Learning algorithms like Naive Bayes, Logistic Regression and Random Forest had been practiced with the dataset. After getting the basic understanding of the ML algorithms, it had been implimented on the Iris Dataset, wherein the confusion matrix, accuracy, precision, recall and fi-score had been calculated and displayed.

```
Naive Bayes on TF-Idf

1    import numpy as np
2    import pandas as pd
3
4    from sklearn.feature_extraction.text import TfidfVectorizer
5
6    from sklearn.naive_bayes import MultinomialNB
7    from sklearn.linear_model import LogisticRegression
8
9    import os
10
11
12   train_df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/train_data_cleaning.csv")
13   test_df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/test_data_cleaning.csv")
14
15   # READ TRAINING DATA AND SEPARATE INTO X AND y
16   X_train = train_df['location']
17   y_train = train_df['text']
18
19   X_test = test_df['location']
20
21   vectorizer = TfidfVectorizer()
22
23   X_train = vectorizer.fit_transform(X_train.values.astype('U'))
24   X_test = vectorizer.transform(X_test.values.astype('U'))
25
26   # MNB CLASSIFICATION
27   mnb = MultinomialNB()
28   mnb.fit(X_train,y_train)
29   mnb_prediction = mnb.predict(X_test)
30
31   mnb_results_TFidf = np.array(list(zip(test_df['id'],mnb_prediction)))
32   mnb_results_TFidf = pd.DataFrame(mnb_results_TFidf, columns=['id', 'text'])
33
34   mnb_results_TFidf.to_csv('/content/drive/MyDrive/Colab Notebooks/mnb_vectorized_TFidf.csv', index = False)
35
36   mnb_results_TFidf
```

Figure 2.4: Naive Bayes on TF-Idf



```
Logistic Regression on Count Vectorizer

1    import numpy as np
2    import pandas as pd
3
4    from sklearn.feature_extraction.text import CountVectorizer
5
6    from sklearn.linear_model import LogisticRegression
7
8    import os
9
10   train_df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/train_data_cleaning.csv")
11   test_df = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/test_data_cleaning.csv")
12
13   # READ TRAINING DATA AND SEPARATE INTO X AND y
14   X_train = train_df['location']
15   y_train = train_df['text']
16
17   X_test = test_df['location']
18
19   vectorizer = CountVectorizer(min_df=4)
20
21   X_train = vectorizer.fit_transform(X_train.values.astype('U'))
22   X_test = vectorizer.transform(X_test.values.astype('U'))
23
24   # LOGISTIC REGRESSION
25   logreg = LogisticRegression()
26   logreg.fit(X_train, y_train)
27   lr_prediction = logreg.predict(X_test)
28
29   lr_results_CV = np.array(list(zip(test_df['id'],lr_prediction)))
30   lr_results_CV = pd.DataFrame(lr_results_CV, columns=['id', 'text'])
31   lr_results_CV.to_csv('lr_vectorized.csv', index = False)
32
33   lr_results_CV
```

| | id | text |
|---|---|---|
| 0 | 0 | The Latest : More Homes Razed by Northern Cal... |
| 1 | 2 | The Latest : More Homes Razed by Northern Cal... |

Figure 2.5: Logistic Regression on CountVectorizer

```
Random forest on Iris

1   import numpy as np
2   import pandas as pd
3
4   from sklearn.model_selection import train_test_split
5   from sklearn import metrics
6
7   from sklearn.feature_extraction.text import CountVectorizer
8
9   from sklearn.naive_bayes import GaussianNB
10  from sklearn.linear_model import LogisticRegression
11  from sklearn.ensemble import RandomForestClassifier
12  from sklearn.metrics import make_scorer, accuracy_score,precision_score
13  from sklearn.metrics import classification_report
14  from sklearn.metrics import confusion_matrix
15  from sklearn.metrics import accuracy_score ,precision_score,recall_score,f1_score
16
17  iris = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/Iris.csv")
18
19  X=iris.iloc[:,0:4].values
20  y=iris["Species"].values
21
22  X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=0)
23
24  random_forest = RandomForestClassifier(n_estimators=100)
25  random_forest.fit(X_train, y_train)
26  Y_prediction = random_forest.predict(X_test)
27  accuracy_rf=round(accuracy_score(y_test,Y_prediction)* 100, 2)
28  acc_random_forest = round(random_forest.score(X_train, y_train) * 100, 2)
29
30
31  cm = confusion_matrix(y_test, Y_prediction)
32  accuracy = accuracy_score(y_test,Y_prediction)
33  precision =precision_score(y_test, Y_prediction,average='micro')
34  recall =  recall_score(y_test, Y_prediction,average='micro')
35  f1 = f1_score(y_test,Y_prediction,average='micro')
36  print('Confusion matrix for Random Forest\n',cm)
37  print('accuracy_random_Forest : %.3f' %accuracy)
38  print('precision_random_Forest : %.3f' %precision)
39  print('recall_random_Forest : %.3f' %recall)
40  print('f1-score_random_Forest : %.3f' %f1)

Confusion matrix for Random Forest
 [[16  0  0]
 [ 0 18  0]
 [ 0  0 11]]
accuracy_random_Forest : 1.000
precision_random_Forest : 1.000
recall_random_Forest : 1.000
f1-score_random_Forest : 1.000
```

Figure 2.6: Random forest on Iris Dataset

## 2.2.6   Task 6

This section deals with the implementation of Convolutional Neural Network on the dataset. For this task, the dataset used is taken from Codalab of the competition on Sentiment Analysis and Homophobia detection of YouTube comments. Out of many datasets to work on, the English dataset had been selected. Firstly, the dataset is imported, and preprocessed. HTML tags, special characters, punctuation are removed. Later, the multiple spaces in the data is cleaned and the entire dataset is rendered in lowercase. Next we removed the unnecessary tags that were in the dataset. We store the preprocessed data in a list. Finally we divide our dataset into train and test sets, followed by tokenizing the lists. Next data embedding is done which converts the textual data into numeric data that is needed for the CNN model. We will create an embedding matrix where each row number will correspond to the index of the word in the corpus the matrix will have two hundred columns where each column will contain the globe word embedding for the words in the corpus. A simple CNN is created with one convolution layer and one pooling layer. A sequential model is created followed by an embedding layer.Convolutional layer is generated with 128 features, the kernel size is 5 and the activation function used here is sigmoid. The feature size is reduced using the pooling layer. Next the model is evaluated and graph is plotted of the results

```
[ ]   1   from keras.models import Model
      2   from keras.layers import Conv1D, Dense, MaxPool1D, Flatten, Input
      3   import numpy as np
      4   model = Sequential()

[ ]   1   from keras.layers.pooling import GlobalMaxPooling1D
      2   embedding_layer = Embedding(vocab_size, 200, weights = [embedding_matrix], input_length = maxlen, trainable = False)
      3   model.add(embedding_layer)
      4
      5   model.add(Conv1D(128, 5, activation = "relu"))
      6   model.add(GlobalMaxPooling1D())
      7   model.add(Dense(1, activation = "sigmoid"))
      8   model.compile(optimizer = "adam", loss = "binary_crossentropy", metrics = ["acc"])

   1   print(model.summary())

Model: "sequential_6"

Layer (type)                  Output Shape            Param #
=================================================================
embedding_7 (Embedding)       (None, 200, 200)        1107600

conv1d_6 (Conv1D)             (None, 196, 128)        128128

global_max_pooling1d_6 (Glo   (None, 128)             0
balMaxPooling1D)

dense_6 (Dense)               (None, 1)               129

=================================================================
Total params: 1,235,857
Trainable params: 128,257
Non-trainable params: 1,107,600

None
```

Figure 2.7: 1D Convolutional Neural Network

### 2.2.7 Task 7

For this task, a dataset containing codemix text had been used. In this the data is preprocessed as the previous tasks. Next in the data we can see there are common short forms that were misspelled or are incorrectly used. For this we create a small dictionary using re and we try to correct the words in our vicinity.

```python
5    #HTML tags
6    sentence = remove_tags(sen)
7
8    #Replacing common shorforms
9    sentence = re.sub(r"i'm", "i am", sentence)
10   sentence = re.sub(r"Plz", "Please", sentence)
11   sentence = re.sub(r"he's", "he is", sentence)
12   sentence = re.sub(r"she's", "she is", sentence)
13   sentence = re.sub(r"that's", "that is", sentence)
14   sentence = re.sub(r"what's", "what is", sentence)
15   sentence = re.sub(r"where's", "where is", sentence)
16   sentence = re.sub(r"\'ll", " will", sentence)
17   sentence = re.sub(r"\'ve", " have", sentence)
18   sentence = re.sub(r"\'re", " are", sentence)
19   sentence = re.sub(r"\'d", " would", sentence)
20   sentence = re.sub(r"\'ve", " have", sentence)
21   sentence = re.sub(r"won't", "will not", sentence)
22   sentence = re.sub(r"don't", "do not", sentence)
23   sentence = re.sub(r"did't", "did not", sentence)
24   sentence = re.sub(r"can't", "can not", sentence)
25   sentence = re.sub(r"it's", "it is", sentence)
26   sentence = re.sub(r"i'ts", "it is", sentence)
27   sentence = re.sub(r"couldn't", "could not", sentence)
28   sentence = re.sub(r"have't", "have not", sentence)
29
30   #Special Characters
31   sentence = re.sub("([^\x00-\x7F])+"," ",sentence)
32
33   #Punctuations
34   sentence = re.sub(r"[,.\"!@#$%^&*(){}?/;`~:<>+=-]", "", sentence)
35
36   #single Char removal
37   sentence = re.sub(r"\s+[a-zA-Z]\s+", ' ', sentence)
38
39   #Multiple spaces
40   sentence = re.sub(r'\s+', ' ', sentence)
41
42   #Lowercase
43   sentence = sentence.lower()
44
45   return sentence
```

Figure 2.8: 1D Convolutional Neural Network

### 2.2.8 Task 8

This task comprises of the work done in the internship. It introduces the data cleaning methods, the use of differnt machine leaning models to process the data and the NLP basics that are applied to the dataset.
First we install the required packages namely textblob, pyspellchecker, emoji, googletrans.
Then we start importing the required packages for the preprocessing of the dataset and also for the visualizing and translation. Next we import the dataset. It is a codemix dataset from kaggle containing tweet. Sentiment

10

analysis is to be done on the dataset.

After importing, we find that the lables are marked 0 fornegative and 4 for the positive reviews. So we create a new column stating the same. We then discard the columns that are not required for the particular NLP task.

After this we start the preprocessing task. First, we remove the emojis, next we compile a code with many contractions and their full forms and apply the same in the dataset. The emoticons, though not used widely nowadays are also taken care of next. After this, we remove the HTML tags, mentions and the stopwords.

We import the required NLTK libraries and then process the data with all the requied functions. The output is a cleaner dataset which we store in a separate file. After this, we correct grammar of the sentences and finally we get a clean text.

We then split the data into train and test lists. We then move forward to feed the data into different classifiers to see their outcomes.

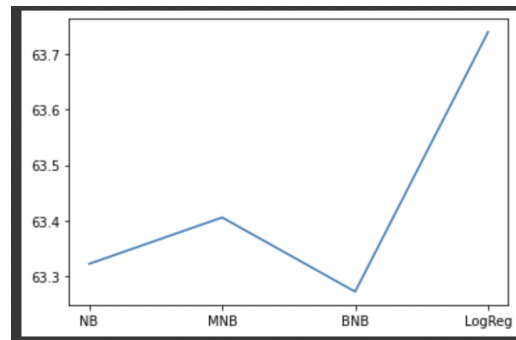| Classifiers | Accuracy |
|---|---|
| Naive Bayes | 63.32277953674388 |
| Multinomail naive bayes | 63.406098983502744 |
| Bernouli naive bayes | 63.27278786868855 |
| Logistic regression | 63.73937677053825 |



Figure 2.9: Classifiers accuracy plot

We notice that the Logistic Regression has the highest accuracy rate. Next we create the confusion matrix of the data (positive and negative tweets). Next we use the google translate library to translate the data and make it monolingual.

Link: **Google Colab Notebook**

# Chapter 3

# Conclusion

I conclude this report by sharing my personal experience with the 2 months internship on Machine Learning and Deep Learning. Firstly, I gained the basic understanding of the working of the python libraries and frameworks. Secondly, I got to work with different datasets and the process of how we can contribute to a competition. Thirdly, I got exposure about the different methods that are to be used in the process of Machine Learning and how those are different from that of Deep Learning. Lastly, I had been practicing the processes and eager to learn more in this domain.

Regarding the future scope of the internship, the models used are not by far perfectly used and there is much more to learn through the libraries used in this internship. The task 7 still contains just the English shortforms and much more data is needed in the real world situations. In task 8 there had been issues with the translation part and not all of the words get correctly translated due to similar words in Hindi and English like "Main", "To", etc.