

News Group Classification

Content :

1. What is the news group classification project ?
 2. Preprocessing
 3. Models learning
 4. Models evaluation
 5. Predictions
 6. Results
 7. Experiment
-

What is the news group classification project ?

The aim of this project is to preprocess a group of articles splitted between classes using NLP techniques and train it using machine learning or deep learning to get a model capable of classifying future articles to its specific news group.

Preprocessing

The data set contain of 20 folders each folder represent a class and inside each folder group of articles in text format

- 1) Remove unwanted headers from each article in the dataset

Define Path and Initialize Dictionary

```
In [ ]: 1 path = "C:\\Users\\Shehab\\Desktop\\NLP\\Dataset"
        2 dic = {}
        3 for folder in os.listdir(path):
        4     dic.update( {folder : []} )
        5
```

Fill the Dictionary with articles and remove first 2 lines

```
In [ ]: 1
        2 for folder in os.listdir(path):
        3     folder_path = os.path.join(path, folder)
        4     for file in os.listdir(folder_path):
        5         file_path = os.path.join(folder_path, file)
        6         with open(file_path, 'r') as f:
        7             contents = f.readlines()
        8             contents = contents[2:]
        9             tmp = dic[folder]
        10            tmp.append(contents)
        11            dic[folder] = tmp
        12
        13
```

- 2) convert the dataset to dataframe so that each row in the data frame consists of 2 columns article and class of the article and encode class labels.

From Dictionary to Data Frames

```
In [ ]: 1 dfs=[]
        2 labels = []
        3 i =0
        4 for cls in dic:
        5     tmp_dic = {"article" : dic[cls] , "group" : [i]*len(dic[cls])}
        6     labels.append(cls)
        7     i+=1
        8     df = pd.DataFrame(tmp_dic)
        9     dfs.append(df)
        10
        11
        12 all_dfs = pd.concat(dfs)
        13 all_dfs.to_csv("All.csv" , index=False)
```

- 3) Tokenize each article in the dataframe using regex Tokenizer and remove stop words

Tokinize and preprocessing

```
In [ ]: 1 all_dfs = pd.read_csv("All.csv")
2 stop = list(stopwords.words('english'))
3 def tokenize(text):
4     text = text.lower()
5     text = text.replace("\\n", "")
6     text = text.replace("\\", "")
7     text = text.replace("'", "")
8     pattern = r"^[a-z]+(-[a-z]+)*\.? $"
9     tokens = word_tokenize(text)
10    tokens = [token for token in tokens if re.match(pattern, token)]
11    tokens = [x for x in tokens if x not in stop]
12
13    return list(tokens)
14
15
16 all_dfs["article"] = all_dfs["article"].apply(tokenize)
17
```

- 4) Apply POS tagging to each article in the data frame then apply lemmatizing to each token and save the data frame in CSV format

Apply POS and Lemma

```
In [ ]: 1 tag_map = {
2     "N": "n",
3     "V": "v",
4     "R": "r",
5     "J": "a"
6 }
7
8 {token: "hello" , tag:"n"}
9 def pos_lemma(tokens):
10     tags=pos_tag(tokens)
11     lemmatizer=WordNetLemmatizer()
12     w_net_lemma=[]
13     for token ,tag in tags:
14         if(tag_map.get(tag[0]) is not None ):
15             w_net_lemma.append(lemmatizer.lemmatize(token,pos=tag_map[tag[0]]))
16         else:
17             w_net_lemma.append(lemmatizer.lemmatize(token,pos='n'))
18
19     return list(w_net_lemma)
20
21 all_dfs["article"] = all_dfs["article"].apply(pos_lemma)
```

Save Data Frame

```
In [ ]: 1 all_dfs.to_csv("All.csv" , index=False)
```

Models Learning

- 1) Load saved CSV dataframe and split it to train and test with the percentage of 80% train and 20 % test

Load Data Frame

```
[13]: 1 all_dfs = pd.read_csv("All.csv")
```

Split Data to Train and Test

```
[14]: 1 train_x, valid_x, train_y, valid_y = train_test_split(all_dfs['article'], all_dfs['group'], test_size=0.2, random_state=42)
```

- 2) Apply TF-IDF on both train and test data frames and save the results in PK1 file and save tfidf_vect to joblib file

Apply TF-IDF on Train and Test Data

```
In [15]: 1 tfidf_vect = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', max_features=5000)
2 tfidf_vect.fit(all_dfs['article'])
3 xtrain_tfidf = tfidf_vect.transform(train_x)
4 xvalid_tfidf = tfidf_vect.transform(valid_x)
```

Save Test data and its Labels

```
In [ ]: 1 joblib.dump(xvalid_tfidf, 'pk1/xvalid_tfidf.pk1')
2 joblib.dump(valid_y, 'pk1/valid_y.pk1')
```

- 3) Declare train template function which take a model data and both train and test data then fit the train data using the given model and show the train/test accuracy plot during training and save model weights to PK1 file

Train Template

```
In [ ]: 1 def All_Models(model,model_name,x_train,y_train,x_test,y_test):
2         #-----model fitting-----
3         history = model.fit(x_train,y_train)
4         #-----Train Plot-----
5         train_sizes, train_scores, test_scores = learning_curve(model, x_train, y_train, cv=5, n_jobs=-1, train_sizes=np.linspace(
6         train_mean = np.mean(train_scores, axis=1)
7         train_std = np.std(train_scores, axis=1)
8         test_mean = np.mean(test_scores, axis=1)
9         test_std = np.std(test_scores, axis=1)
10        plt.figure(figsize=(8, 6))
11        #plt.plot(train_sizes, train_mean, 'o-', color='b', label='Training score')
12        plt.plot(train_sizes, test_mean, 'o-', color='r', label='Cross-validation score')
13        #plt.fill_between(train_sizes, train_mean - train_std, train_mean + train_std, alpha=0.2, color='b')
14        plt.fill_between(train_sizes, test_mean - test_std, test_mean + test_std, alpha=0.2, color='r')
15        plt.xlabel('Number of training samples')
16        plt.ylabel('Score')
17        plt.title('Learning Curve')
18        plt.legend(loc='best')
19        plt.show()
20        #-----save model-----
21        joblib.dump(model, f'pk1/{model_name}.pk1')
22        #####
23
```

4) Pass models to the train template function

Naive Bayes

```
In [ ]: 1 model = naive_bayes.MultinomialNB(alpha=0.2)
2 All_Models(model,"Naive Bayes Multinomial",xtrain_tfidf,train_y,xvalid_tfidf,valid_y)
```

Logistic Regression

```
In [ ]: 1 model = linear_model.LogisticRegression(max_iter=10000000)
2 All_Models(model,"Logistic Regression",xtrain_tfidf,train_y,xvalid_tfidf,valid_y)
```

Support Vector Machine

```
In [ ]: 1 model = svm.SVC()
2 All_Models(model,"Support Vector Classifier",xtrain_tfidf,train_y,xvalid_tfidf,valid_y)
```

Random Forest

```
In [ ]: 1 model = RandomForestClassifier(n_estimators=100, random_state=42)
2 All_Models(model,"Random Forest Classifier",xtrain_tfidf,train_y,xvalid_tfidf,valid_y)
```

- 5) Train a neural network that consists of 7 layers of dropout to reduce overfitting and dense layers with relu function and classification layer with 20 class and softmax activation function. Compile the network using adam optimizer and train it then save the weights to H5 file

Neural Network

```
In [18]: 1 model = Sequential()
2 model.add(Dense(256, activation='relu', input_dim=xtrain_tfidf.shape[1]))
3 model.add(Dropout(0.5))
4 model.add(Dense(128, activation='relu'))
5 model.add(Dropout(0.5))
6 model.add(Dense(64, activation='relu'))
7 model.add(Dropout(0.5))
8 model.add(Dense(20, activation='softmax'))
9 model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
10 x = xtrain_tfidf.toarray()
11 x_valid = xvalid_tfidf.toarray()
12 model.fit(x, train_y, epochs=10, batch_size=32, verbose=1)
13 model.save("pk1/NN_model.h5")
14
15
```

Model evaluation

- 1) Load tf-idf test data and its labels from the PK1 files then declare an accuracy calc function that uses the passed model weights and information to calculate accuracy and plot confusion matrix using test data.

Load Test data and its Labels

```
In [3]: 1 x_test = joblib.load('pk1/xvalid_tfidf.pk1')
2 y_test = joblib.load('pk1/valid_y.pk1')
```

Accuracy Calculation Template

```
In [4]: 1 def accuracy_calc(model, model_name):
2     predictions = model.predict(x_test)
3     print(f"Accuracy For {model_name} is = {metrics.accuracy_score(predictions, y_test)}")
4     #-----Confusion Matrix-----
5     cm = confusion_matrix(y_test, predictions)
6     cm_percent = cm / cm.sum(axis=1, keepdims=True) * 100
7     class_names = y_test.unique()
8     plt.figure(figsize=(10, 10))
9     sns.heatmap(cm_percent, annot=True, fmt=".1f", cmap="Blues")
10    plt.xlabel("Predicted Labels")
11    plt.ylabel("True Labels")
12    plt.title("Confusion Matrix Heatmap")
13    plt.xticks(np.arange(len(class_names)) + 0.5, class_names)
14    plt.yticks(np.arange(len(class_names)) + 0.5, class_names)
15    plt.show()
```

- 2) Load Models PK1 files and pass them to accuracy calc function to calculate accuracy and plot the confusion matrix

Naive Bayes Accuracy ¶

```
1 naive_bayes_model = joblib.load('pk1/Naive Bayes Multinomial.pk1')
2 accuracy_calc(naive_bayes_model , "Naive Bayes Multinomial" )
3 #Accuracy For Naive Bayes Multinomial is = 0.8661710037174721
```

Accuracy For Naive Bayes Multinomial is = 0.8661710037174721

Logistic Regression Accuracy

```
1 logistic_regression_model = joblib.load('pk1/Logistic Regression.pk1')
2 accuracy_calc(logistic_regression_model , "Logistic Regression" )
3 #Accuracy For Logistic Regression is = 0.8629845990440786
```

Accuracy For Logistic Regression is = 0.8629845990440786

Support Vector Machine Accuracy

```
1 svm_model = joblib.load('pk1/Support Vector Classifier.pk1')
2 accuracy_calc(svm_model , "Support Vector Classifier" )
3 #Accuracy For Support Vector Classifier is = 0.8757302177376527
```

Accuracy For Support Vector Classifier is = 0.8757302177376527

Random Forest Accuracy

```
: 1 RNF_model = joblib.load('pk1/Random Forest Classifier.pk1')
2 accuracy_calc(RNF_model , "Random Forest Classifier" )
3 #Accuracy For Random Forest Classifier is = 0.8066914498141264
```

Accuracy For Random Forest Classifier is = 0.8066914498141264

3) Use neural network weights saved in H5 file to evaluate the neural network on test data

Neural network Accuracy

```
1 NN_model = load_model("pk1/NN_model.h5")
2 x_test_reshaped = x_test.toarray()
3 loss, accuracy = NN_model.evaluate(x_test_reshaped, y_test)
4 print(f"Loss is {loss}")
5 print(f"accuracy is {accuracy}")
6 #accuracy is 0.8903345465660095 (BEST)
```

Predictions :

- 1) Save any article to predict.txt file
- 2) Load the predict.txt file and save its content to string

```
1 def load_file_to_string(file_path):
2     with open(file_path, 'r') as file:
3         text = file.read()
4     return text
5
6 file_path = 'predict.txt'
7 file_contents = load_file_to_string(file_path)
8
```

- 3) Tokenize the text and remove stop words

```
: 1 stop = list(stopwords.words('english'))
2 def tokenize(text):
3     text = text.lower()
4     text = text.replace("\\n", "")
5     text = text.replace("\\", "")
6     text = text.replace("'", "")
7     pattern = r"^[a-z]+(-[a-z]+)*\.$"
8     tokens = word_tokenize(text)
9     tokens = [token for token in tokens if re.match(pattern, token)]
10    tokens = [x for x in tokens if x not in stop]
11
12    return list(tokens)
13
14 text_token = tokenize(file_contents)
```

- 4) apply pos tagging and lemmatizing


```

tag_map = {
    "N": "n",
    "V": "v",
    "R": "r",
    "J": "a"
}

def pos_lemma(tokens):
    tags=pos_tag(tokens)
    lemmatizer=WordNetLemmatizer()
    w_net_lemma=[]
    for token ,tag in tags:
        if(tag_map.get(tag[0]) is not None ):
            w_net_lemma.append(lemmatizer.lemmatize(token,pos=tag_map[tag[0]]))
        else:
            w_net_lemma.append(lemmatizer.lemmatize(token,pos='n'))

    return list(w_net_lemma)

text_lemma = pos_lemma(text_token)

```

- 5) Load the tf-idf of the dataset and transform the text based on it

```

1 tfidf_vect = joblib.load("pk1/general_tfidf.joblib")
2 text_tfidf = tfidf_vect.transform(text_lemma)
3

```

- 6) Initialize dictionary that have every class name and encoded number from the preprocessing phase

```

1 class_dict = {
2     0: 'alt.atheism',
3     1: 'comp.graphics',
4     2: 'comp.os.ms-windows.misc',
5     3: 'comp.sys.ibm.pc.hardware',
6     4: 'comp.sys.mac.hardware',
7     5: 'comp.windows.x',
8     6: 'misc.forsale',
9     7: 'rec.autos',
10    8: 'rec.motorcycles',
11    9: 'rec.sport.baseball',
12    10: 'rec.sport.hockey',
13    11: 'sci.crypt',
14    12: 'sci.electronics',
15    13: 'sci.med',
16    14: 'sci.space',
17    15: 'soc.religion.christian',
18    16: 'talk.politics.guns',
19    17: 'talk.politics.mideast',
20    18: 'talk.politics.misc',
21    19: 'talk.religion.misc'
22 }

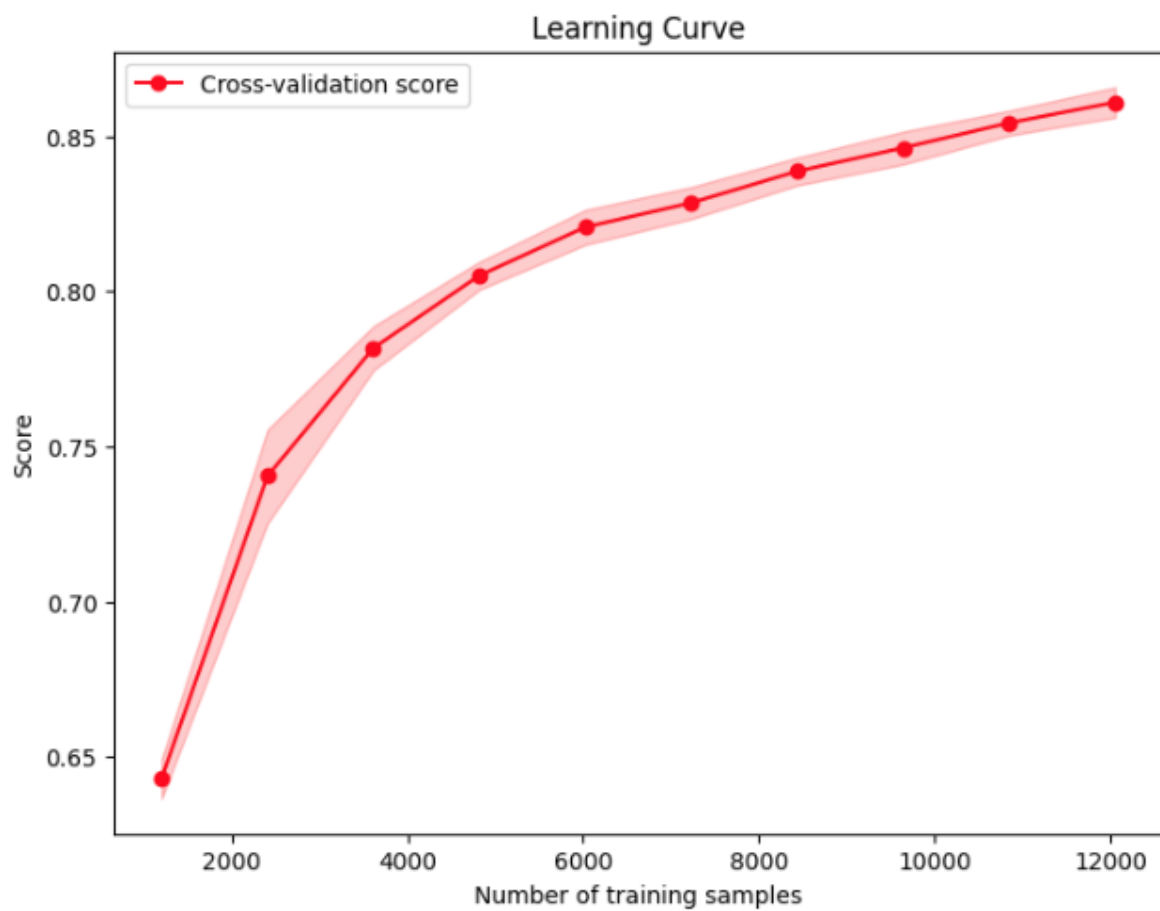
```

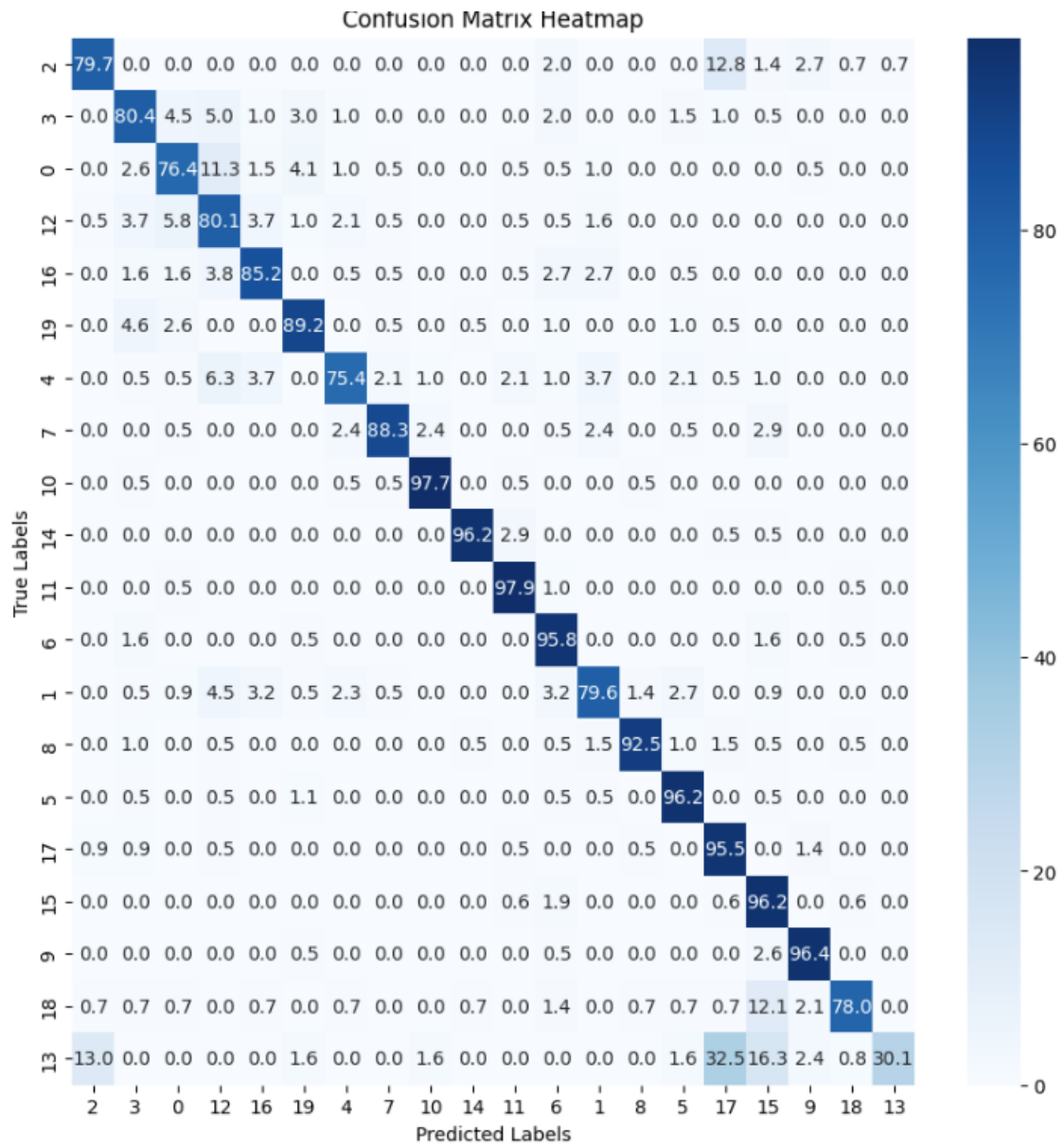
7) Load NN weights and predict the class based on it

```
NN_model = load_model("pk1/NN_model.h5")
text_tfidf_resaped = text_tfidf.toarray()
pred = NN_model.predict(text_tfidf_resaped)
print(f"Predictions for NN is {pred}")
```

Results

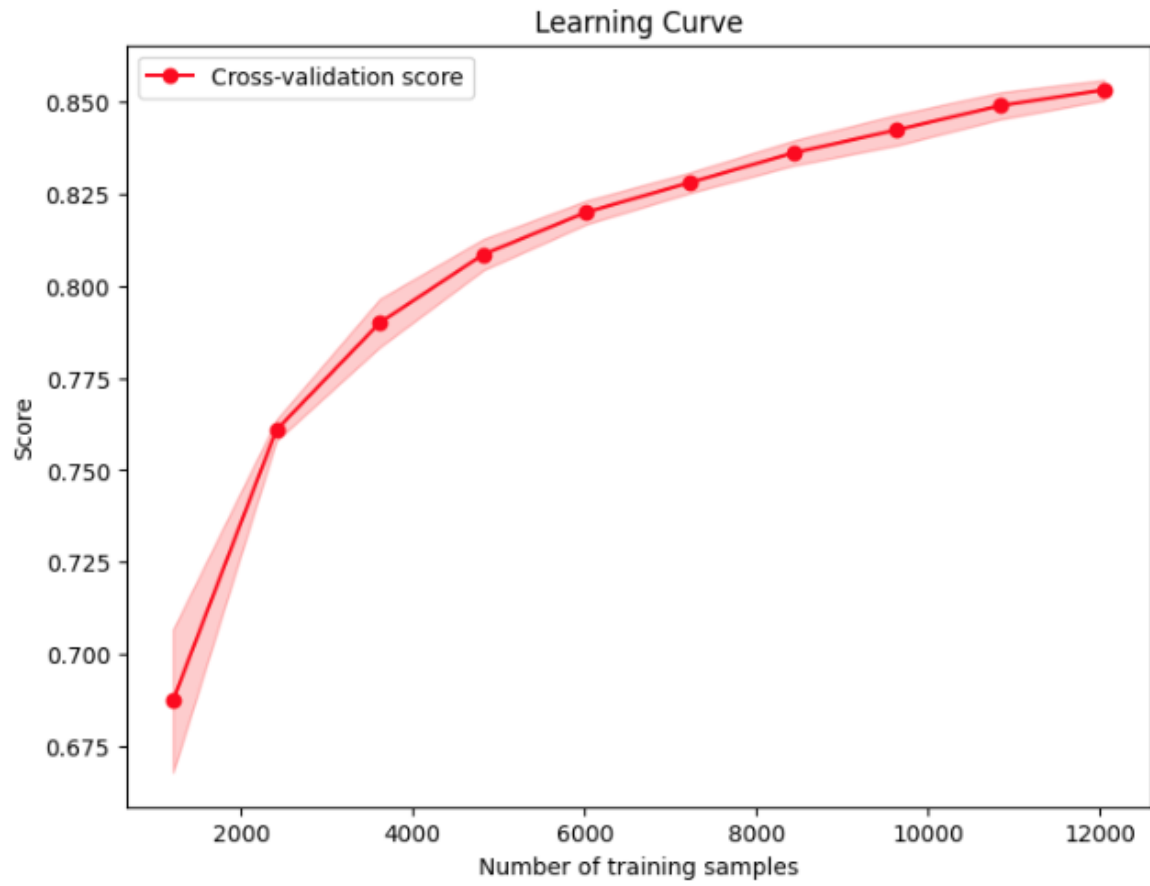
- 1) Naive bayes
Accuracy is 86%

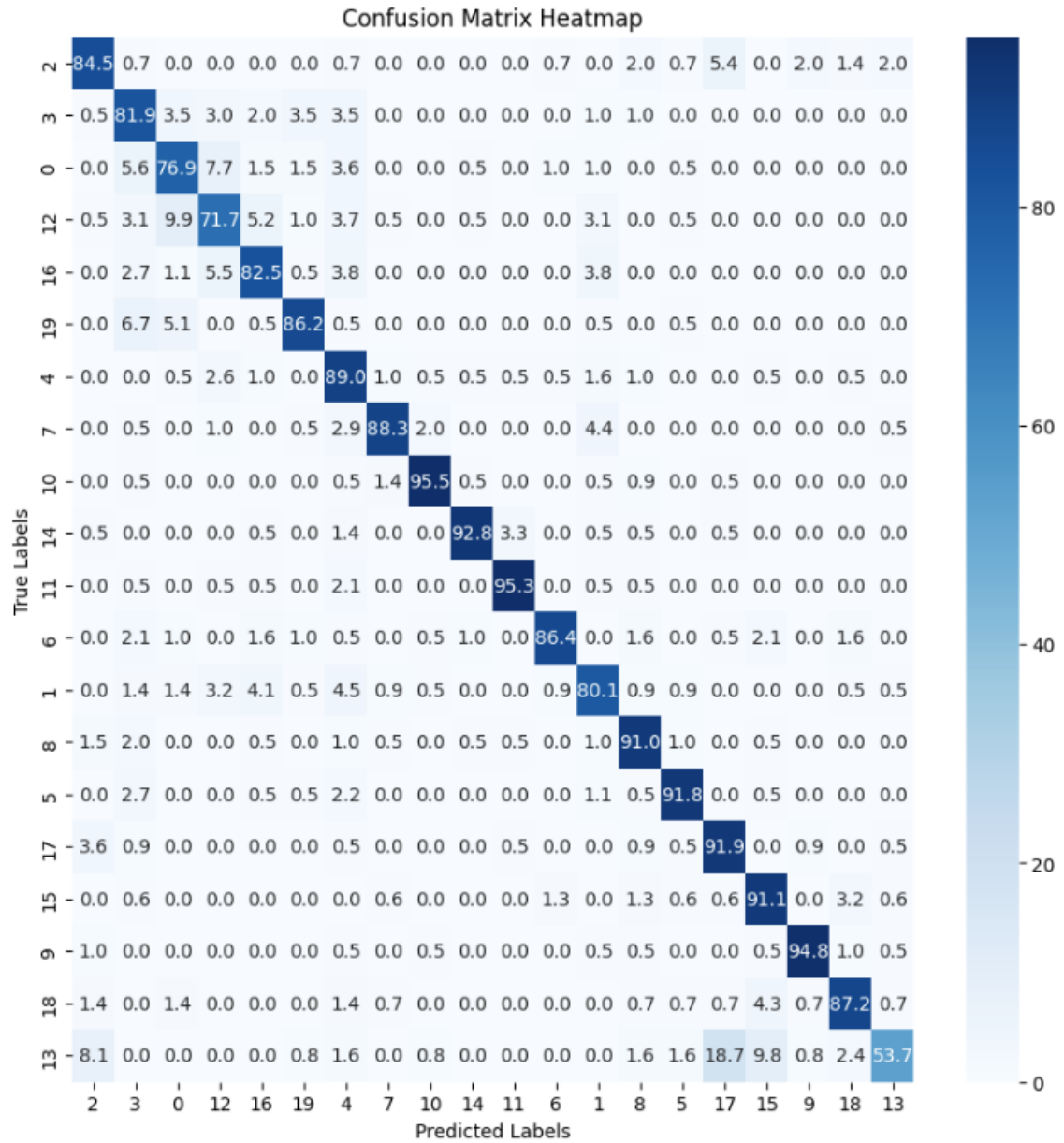




2) Logistic Regression

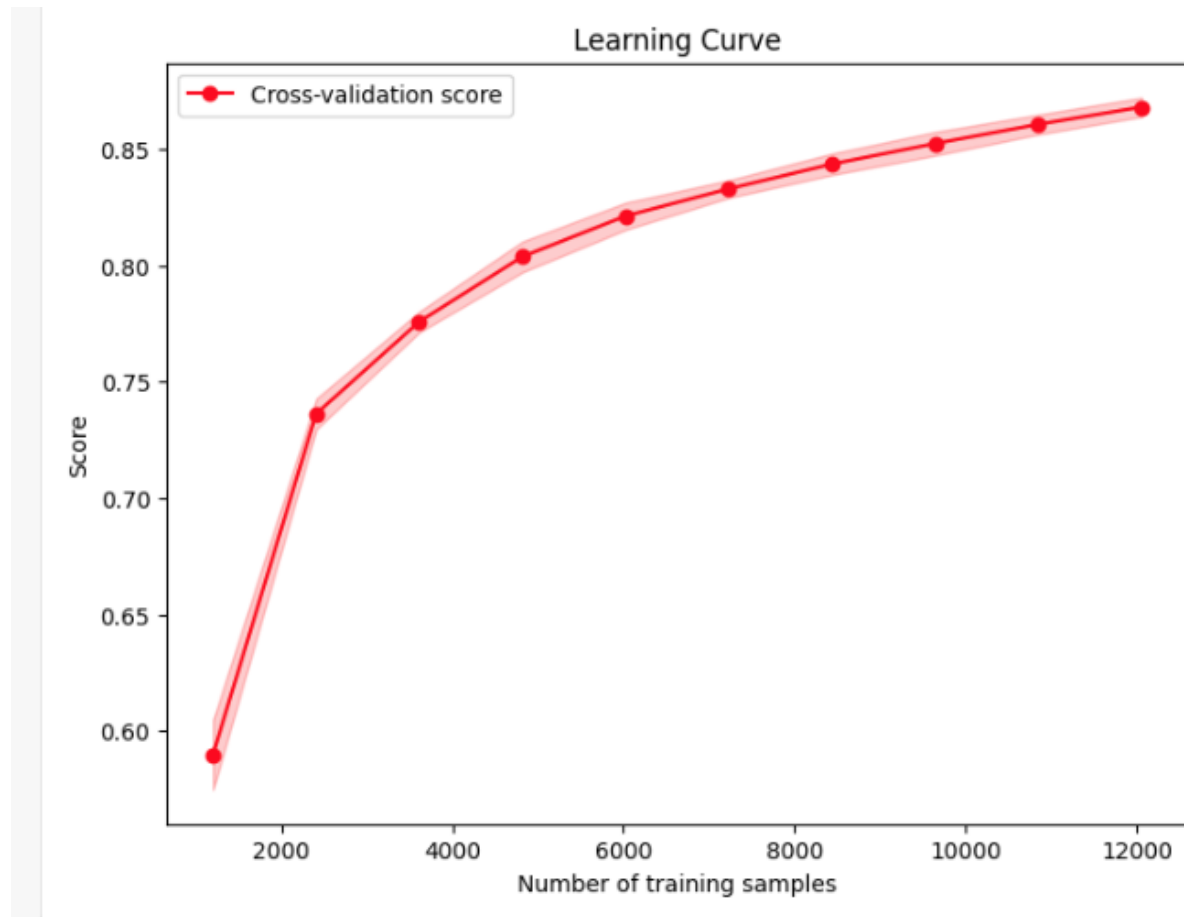
Accuracy is 86%

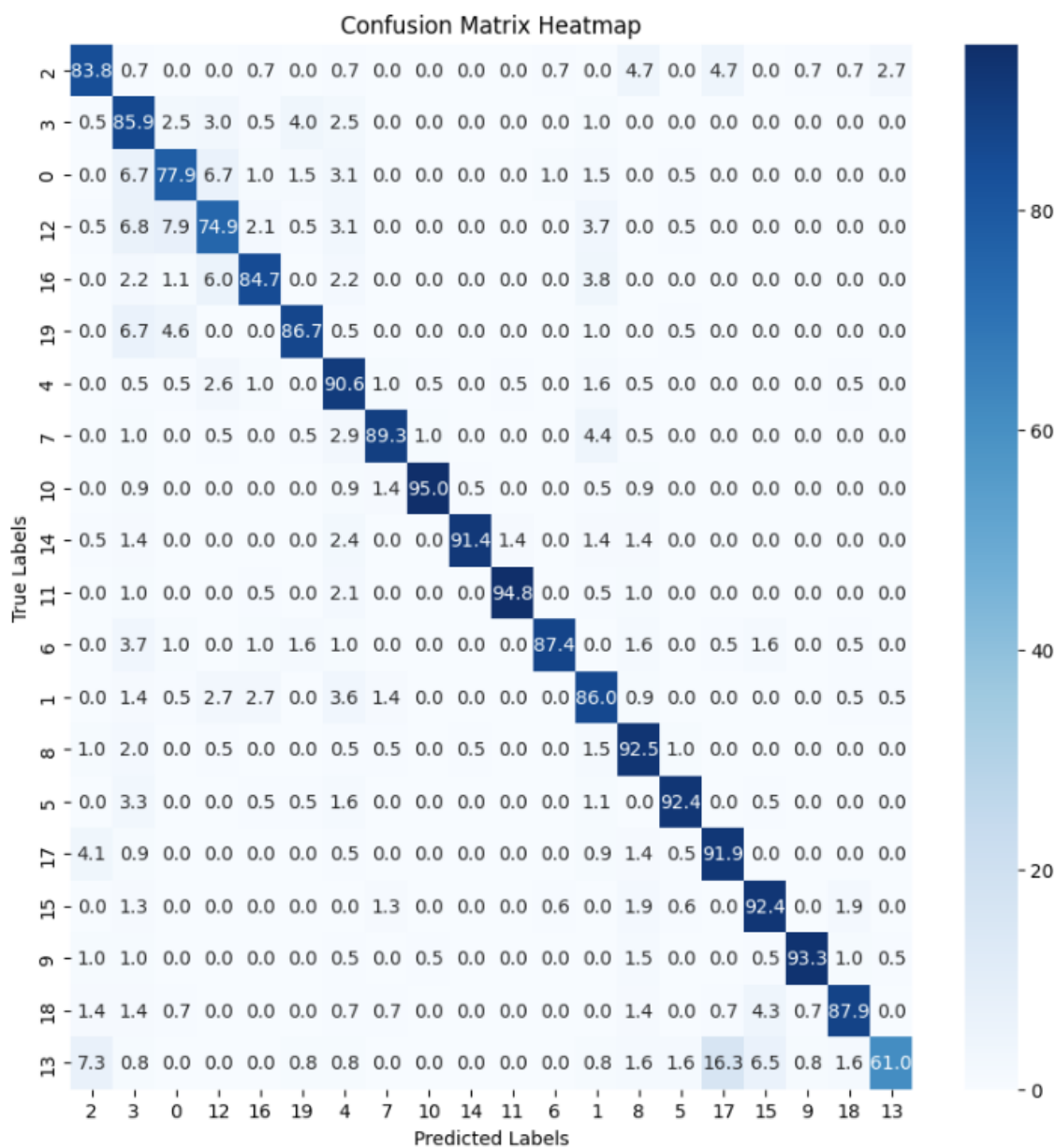




3) Support Vector Machine

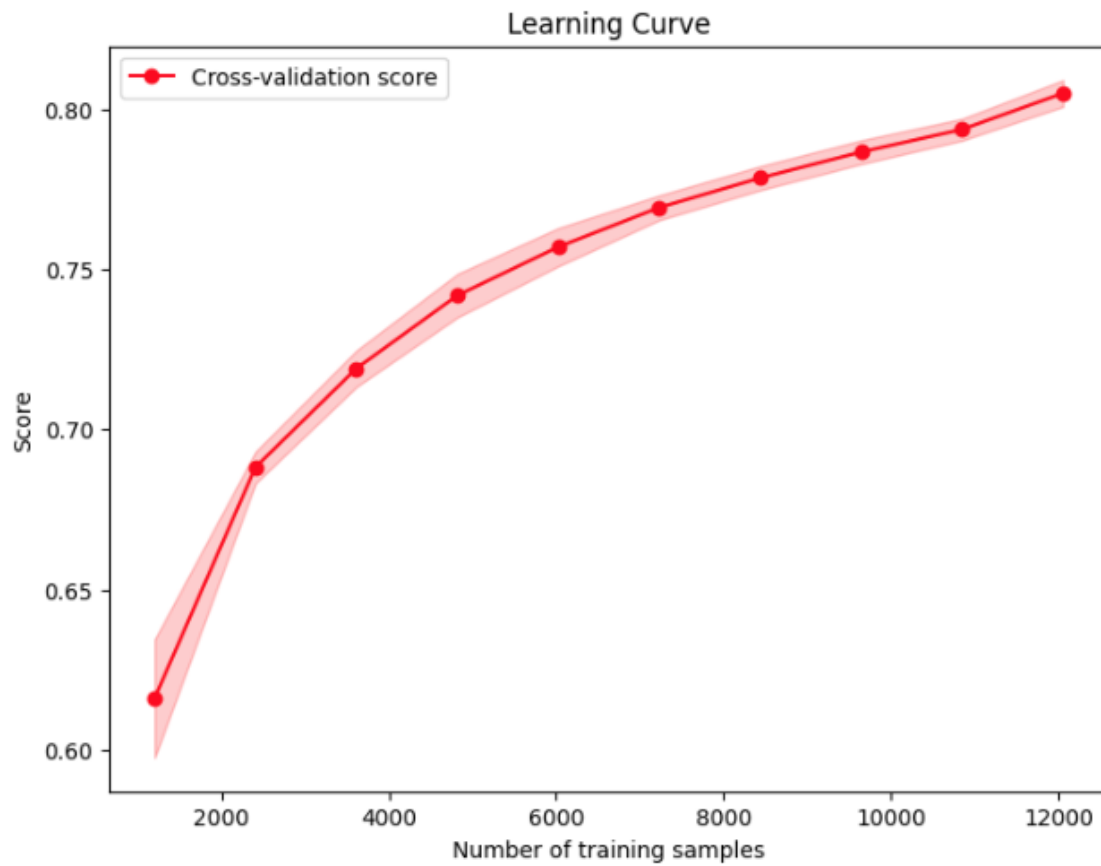
Accuracy is 87%



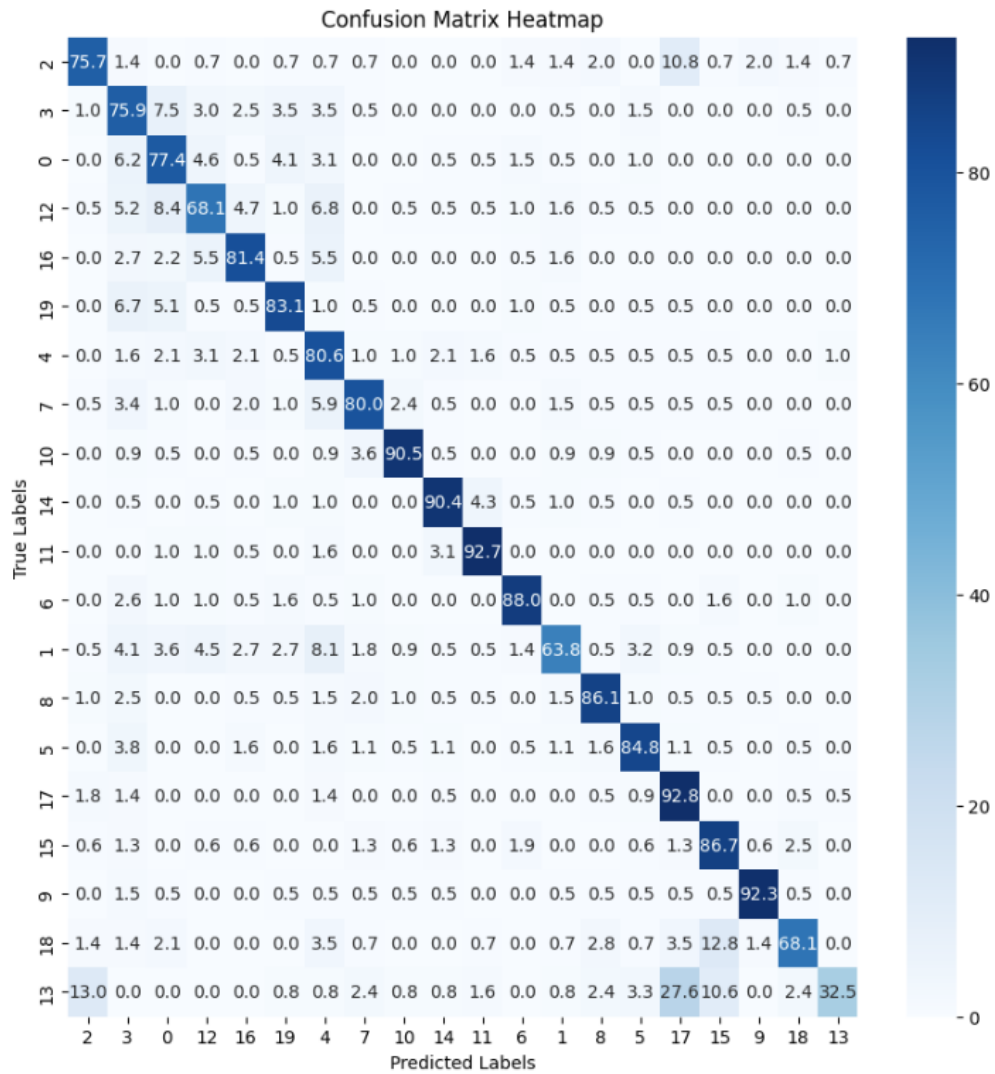


4) Random Forest

Accuracy is 80%



Accuracy For Random Forest Classifier is = 0.8066914498141264

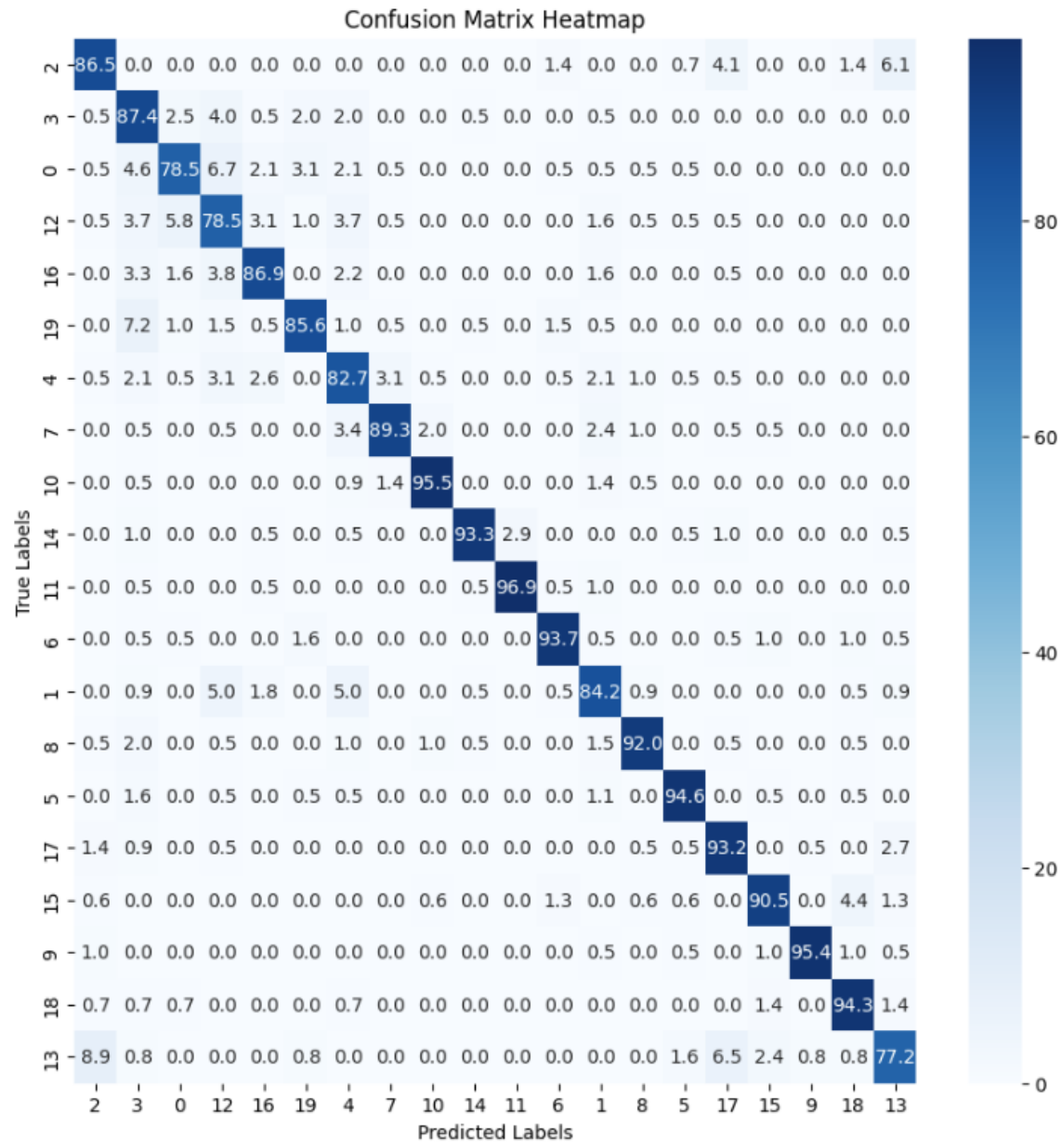


5) Neural Network

Train Accuracy is 98%

Test Accuracy is 89% (BEST Accuracy)

118/118 [=====] - 2s 5ms/step - loss: 0.6431 - accuracy: 0.8903
Loss is 0.6430635452270508
accuracy is 0.8903345465660095



Experiment

Using text = “computer consists of electrical components. I use windows 10 because i don't like windows 11”

The predictions from each model is :

```
1/1 [=====] - 1s 544ms/step  
Predictions for NN is comp.sys.ibm.pc.hardware
```