

# experimentdata

May 1, 2021

## 1 include Libraries

```
[55]: import os
from sklearn import datasets
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.utils import shuffle
import numpy as np
from sklearn.preprocessing import LabelEncoder
from sklearn.pipeline import Pipeline
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import LinearSVC
from sklearn.naive_bayes import BernoulliNB
from sklearn.model_selection import train_test_split
from sklearn.pipeline import Pipeline
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn import tree
from sklearn.neighbors import KNeighborsClassifier
import seaborn as sns
from sklearn.metrics import confusion_matrix
from sklearn.linear_model import LinearRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.feature_extraction.text import TfidfVectorizer
import warnings

import pickle

warnings.filterwarnings("ignore")
plt.style.use('seaborn-darkgrid')
```

```
[56]: # display font setting
font = {'family' : 'normal',
```

```
        'weight' : 'bold',  
        'size'    : 18}  
plt.rc('font', **font)
```

## 2 Read Data

```
[57]: xls = pd.ExcelFile('HS-RU-20.xlsx')  
NeutralHostile = pd.read_excel(xls, 'NeutralHostile')  
HateOffensive = pd.read_excel(xls, 'HateOffensive')
```

## 3 Preprocess data

```
[58]: NeutralHostile.shape
```

```
[58]: (5000, 22)
```

```
[59]: NeutralHostile.columns
```

```
[59]: Index(['Sentence', 'label', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4',  
          'Unnamed: 5', 'Unnamed: 6', 'Unnamed: 7', 'Unnamed: 8', 'Unnamed: 9',  
          'Unnamed: 10', 'Unnamed: 11', 'Unnamed: 12', 'Unnamed: 13',  
          'Unnamed: 14', 'Unnamed: 15', 'Unnamed: 16', 'Unnamed: 17',  
          'Unnamed: 18', 'Unnamed: 19', 'Unnamed: 20', 'Unnamed: 21'],  
          dtype='object')
```

```
[60]: df1 = NeutralHostile.iloc[:, :2].copy()
```

```
[61]: df1.shape
```

```
[61]: (5000, 2)
```

```
[62]: df1['label'] = df1['label'].str.replace('N ', 'N')
```

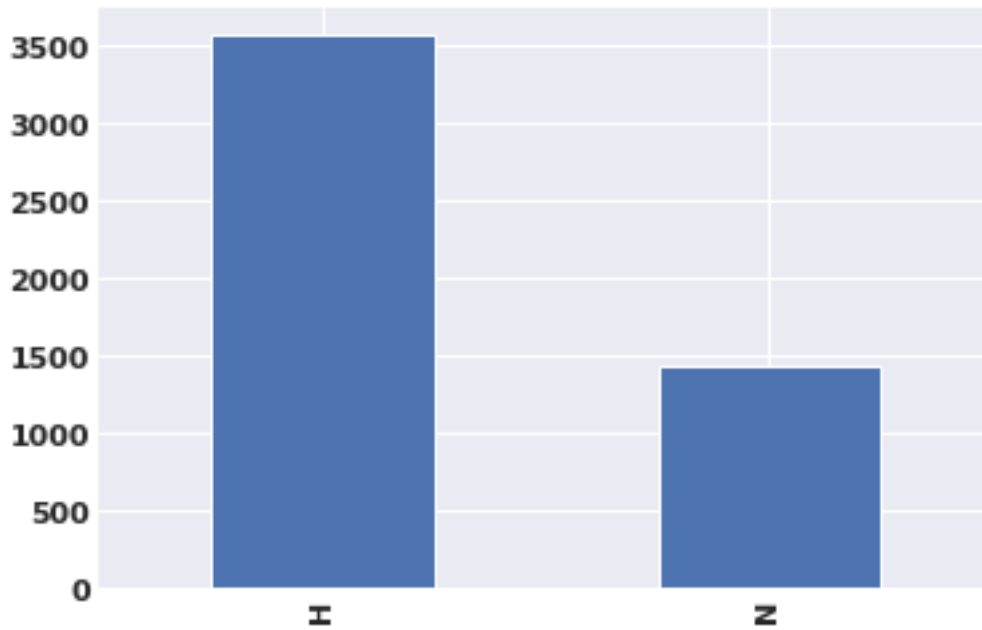
```
[63]: del NeutralHostile
```

```
[64]: output = df1.label.value_counts()
```

```
[65]: output.plot(kind='bar')
```

```
[65]: <AxesSubplot:>
```

findfont: Font family ['normal'] not found. Falling back to DejaVu Sans.



```
[66]: output
```

```
[66]: H      3574
      N      1426
      Name: label, dtype: int64
```

## 4 Word Trigram

```
[67]: vectorizer = TfidfVectorizer(ngram_range=(3,3), analyzer='word', max_features=1000)
```

```
[68]: df1values = df1['Sentence'].values
```

```
[69]: x = vectorizer.fit_transform(df1values)
```

```
[70]: df1v1 = pd.DataFrame(x.toarray(), columns=vectorizer.get_feature_names())
```

```
[71]: df1v1.shape
```

```
[71]: (5000, 1000)
```

## 5 Concatenate Label with prepared data

```
[72]: df1final = pd.concat([df1v1,df1['label']],axis=1)
```

```
[73]: df1final.shape
```

```
[73]: (5000, 1001)
```

```
[74]: df1final.to_excel("ProcessedNeutralHostile.xlsx",index=False)
```

## 6 Label Encodeing

```
[75]: labelEncoder = LabelEncoder()
```

```
[76]: df1encoded_final = df1final.copy()
```

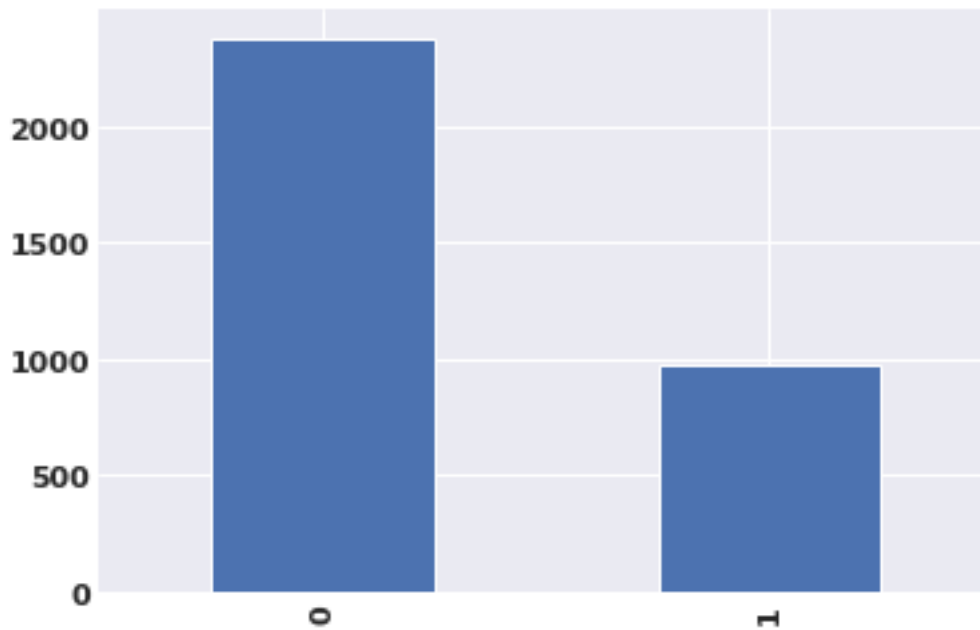
```
[77]: df1encoded_final['label'] = labelEncoder.fit_transform(df1final['label'])
```

## 7 Split dataset

```
[78]: df1encoded_final_train, df1encoded_final_test =  
↳train_test_split(shuffle(df1encoded_final),test_size=.33)
```

```
[79]: trainclasses = df1encoded_final_train.label.value_counts()  
trainclasses.plot(kind='bar')
```

```
[79]: <AxesSubplot:>
```



```
[80]: trainclasses
```

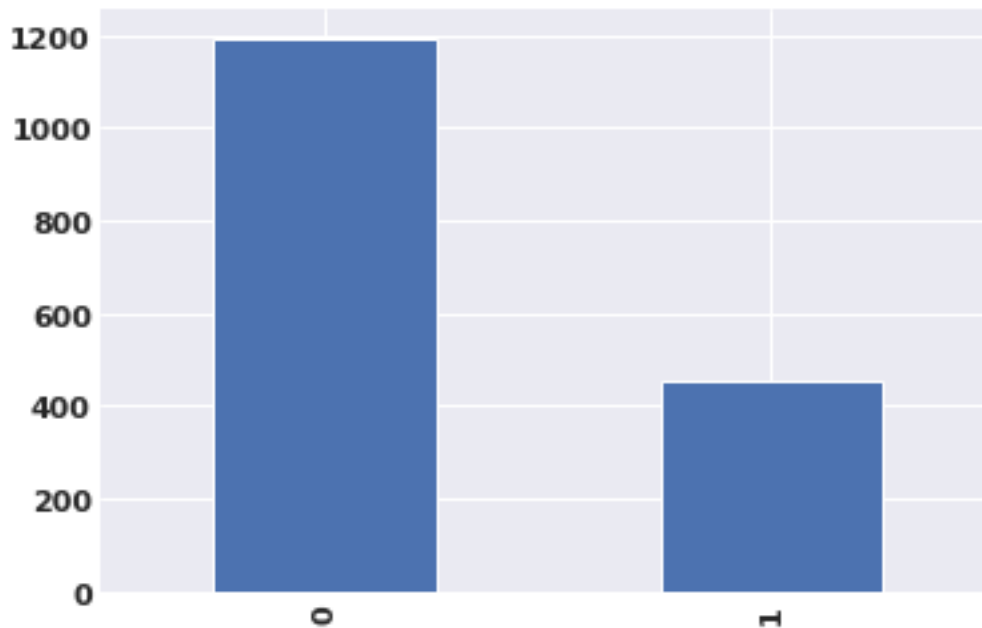
```
[80]: 0    2380  
      1     970  
      Name: label, dtype: int64
```

```
[81]: test_classes = df1encoded_final_test.label.value_counts()  
      test_classes
```

```
[81]: 0    1194  
      1     456  
      Name: label, dtype: int64
```

```
[82]: test_classes.plot(kind='bar')
```

```
[82]: <AxesSubplot:>
```



```
[83]: trainX = df1encoded_final_train.values[:, :-1]
      trainY = df1encoded_final_train["label"].values
      testX = df1encoded_final_test.values[:, :-1]
      testY = df1encoded_final_test["label"].values
```

```
[84]: train = df1final
```

```
[85]: train['label'].unique()
```

```
[85]: array(['H', 'N'], dtype=object)
```

```
[ ]:
```

```
[86]: # confusion matrix
def print_confusion_matrix(name, confusion_matrix, class_names, figsize = (8,5),
    ↪fontsize=14):
    """Prints a confusion matrix, as returned by sklearn.metrics.
    ↪confusion_matrix, as a heatmap.

    Arguments
    -----
    confusion_matrix: numpy.ndarray
        The numpy.ndarray object returned from a call to sklearn.metrics.
    ↪confusion_matrix.
        Similarly constructed ndarrays can also be used.
    class_names: list
```

An ordered list of class names, in the order they index the given  
 ↳confusion matrix.

figsize: tuple  
 A 2-long tuple, the first value determining the horizontal size of the  
 ↳ouputted figure,  
 the second determining the vertical size. Defaults to (10,7).

fontsize: int  
 Font size for axes labels. Defaults to 14.

Returns  
 -----  
 matplotlib.figure.Figure  
 The resulting confusion matrix figure

```

"""
df_cm = pd.DataFrame(
    confusion_matrix, index=class_names, columns=class_names,
)
fig = plt.figure(figsize=figsize)
try:
    heatmap = sns.heatmap(df_cm, annot=True, fmt="d")
except ValueError:
    raise ValueError("Confusion matrix values must be integers.")
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0,
↳ha='right', fontsize=fontsize)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=45,
↳ha='right', fontsize=fontsize)
plt.ylabel('True label', fontsize=22)
plt.xlabel('Predicted label', fontsize=22)
plt.savefig(name+'.jpeg',dpi=600,quality=100,format='jpeg',pad_inches=0.1,
    transparent=True, bbox_inches='tight')

```

```

[87]: bernoulliNB =BernoulliNB()
bernoulliNB.fit(trainX,trainY)
BNBScores = cross_val_score(bernoulliNB,trainX ,trainY,
↳cv=10,scoring='accuracy')
BNBScore = bernoulliNB.score(testX,testY)
print("Accuracy Score : %f"%(BNBScore))
print("Accuracy Average Score of K Fold : %f"%(BNBScores.mean()))

bnb_perdiction = bernoulliNB.predict(testX)

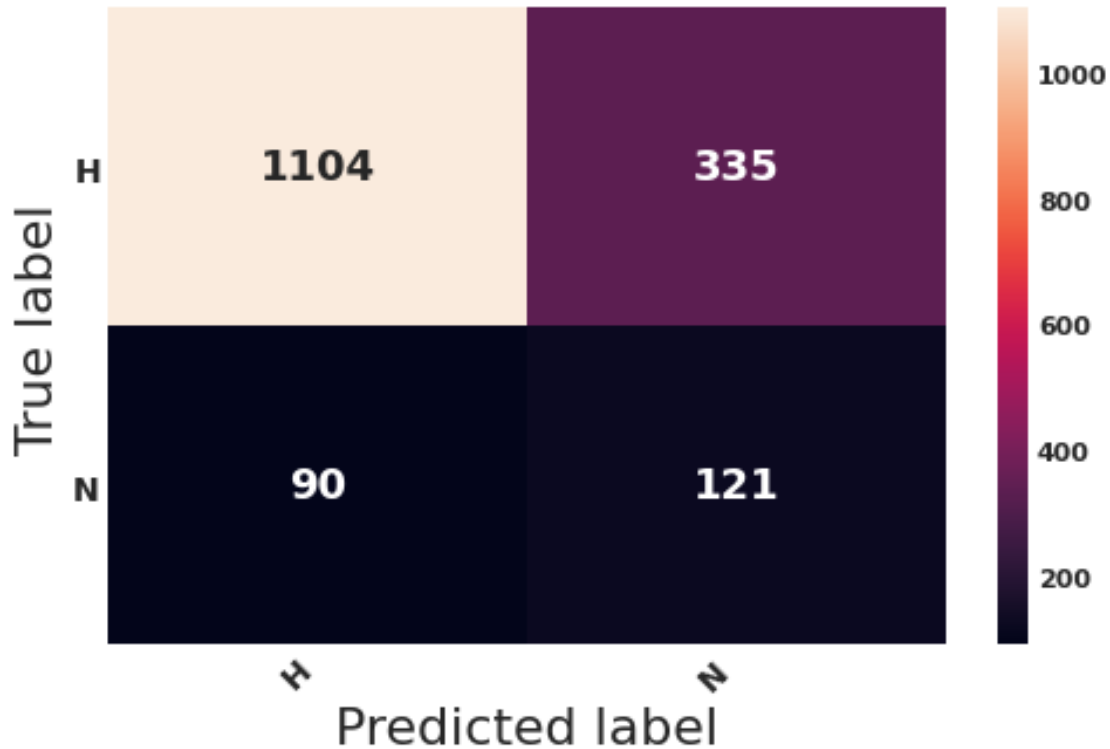
bnb_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(bnb_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
#bnb_con_matrix

```

```
print_confusion_matrix('BernoulliNB',bnb_con_matrix,train['label'].unique())
```

Accuracy Score : 0.742424

Accuracy Average Score of K Fold : 0.727164



```
[88]: linearSVC=LinearSVC()
linearSVC.fit(trainX,trainY)
LSVCScores = cross_val_score(linearSVC,trainX ,trainY, cv=10,scoring='accuracy')
LSVCScore = linearSVC.score(testX,testY)
print("Accuracy Score : %f"%(LSVCScore))
print("Accuracy Average Score of K Fold : %f"%(LSVCScores.mean()))

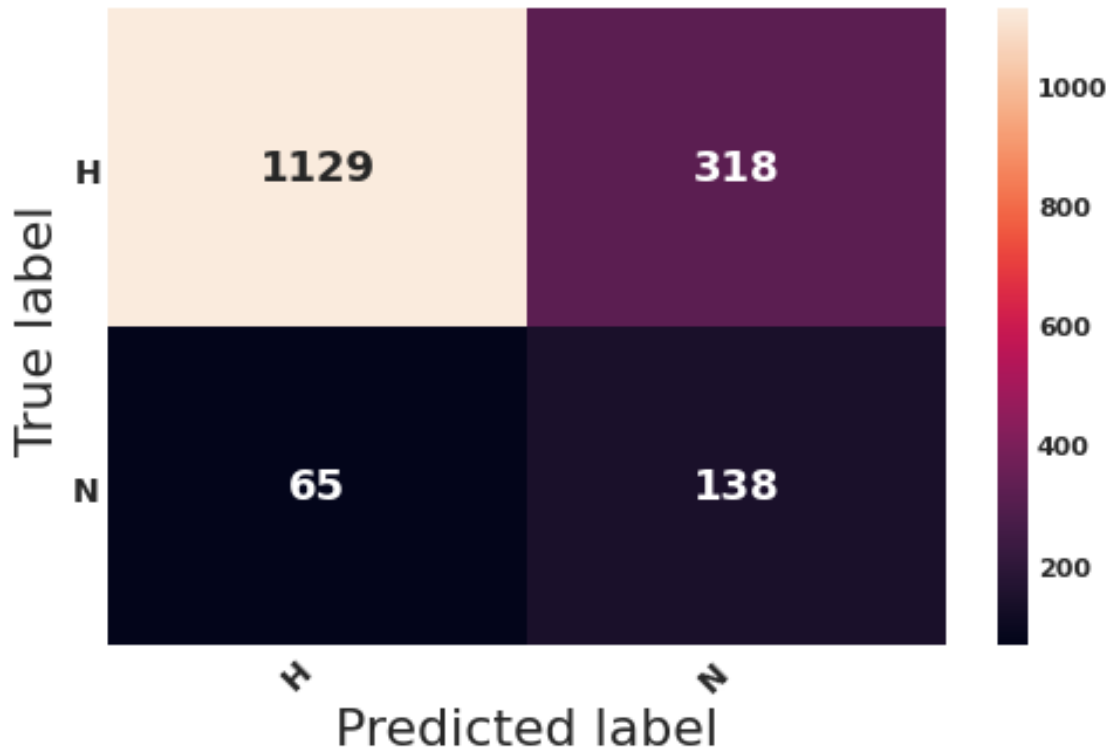
lsvc_perdiction = linearSVC.predict(testX)

lsvc_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(lsvc_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('LinearSVC',lsvc_con_matrix,train['label'].unique())
```

Accuracy Score : 0.767879

Accuracy Average Score of K Fold : 0.749254





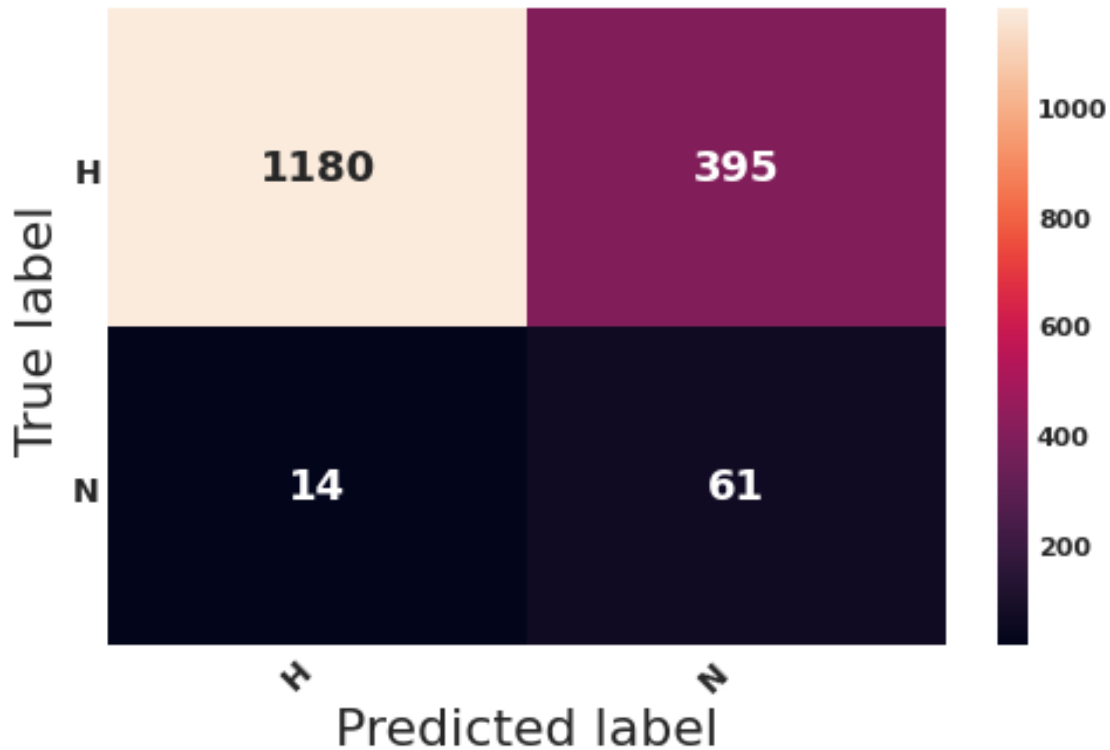
```
[89]: logisticRegression = LogisticRegression();
logisticRegression.fit(trainX,trainY)
LRScores = cross_val_score(logisticRegression,trainX ,trainY,
    ↪cv=10,scoring='accuracy')
LRScore = logisticRegression.score(testX,testY)
print("Accuracy Score : %f"%(LRScore))
print("Accuracy Average Score of K Fold : %f"%(LRScores.mean()))

logisticRegression_prediction = logisticRegression.predict(testX)

LR_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(logisticRegression_prediction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('logisticRegression',LR_con_matrix,train['label'].
    ↪unique())
```

Accuracy Score : 0.752121

Accuracy Average Score of K Fold : 0.733134



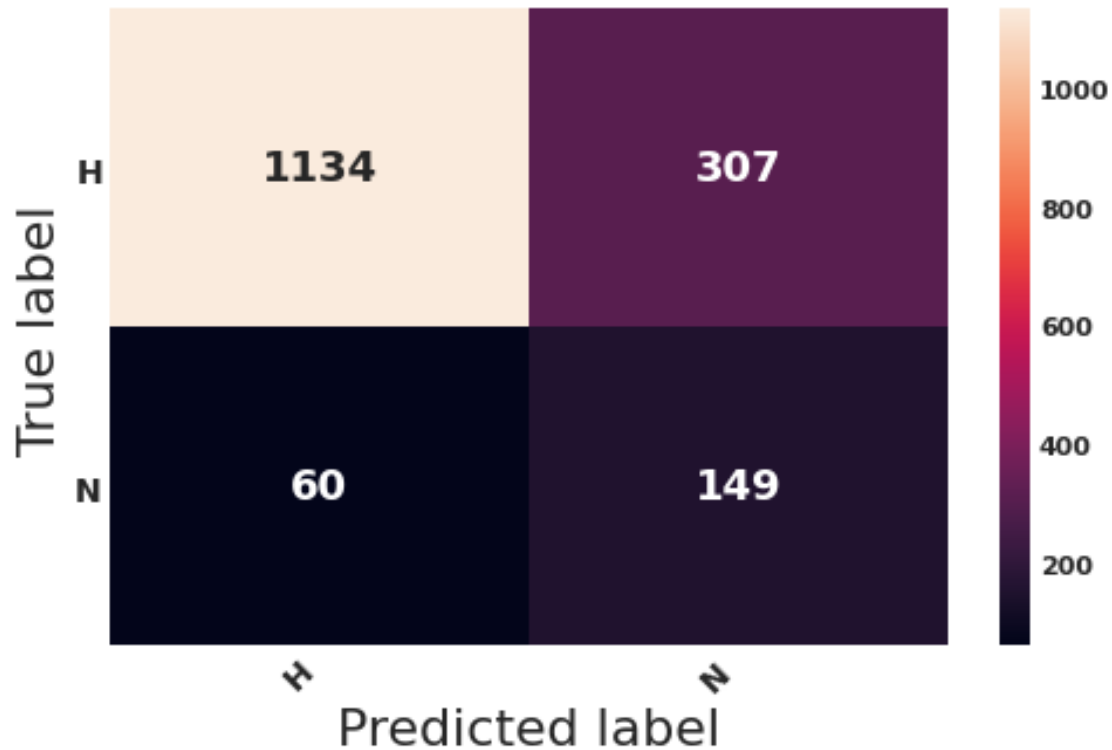
```
[90]: decisionTreeClassifier = tree.DecisionTreeClassifier()
decisionTreeClassifier.fit(trainX,trainY)
DTCScores = cross_val_score(decisionTreeClassifier,trainX ,trainY,
    ↪cv=10,scoring='accuracy')
DTCScore = decisionTreeClassifier.score(testX,testY)
print("Accuracy Score : %f"%(DTCScore))
print("Accuracy Average Score of K Fold : %f"%(DTCScores.mean()))

DTCS_perdiction = decisionTreeClassifier.predict(testX)

DTC_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(DTCS_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('DecisionTreeClassifier',DTC_con_matrix,train['label'].
    ↪unique())
```

Accuracy Score : 0.777576

Accuracy Average Score of K Fold : 0.755522



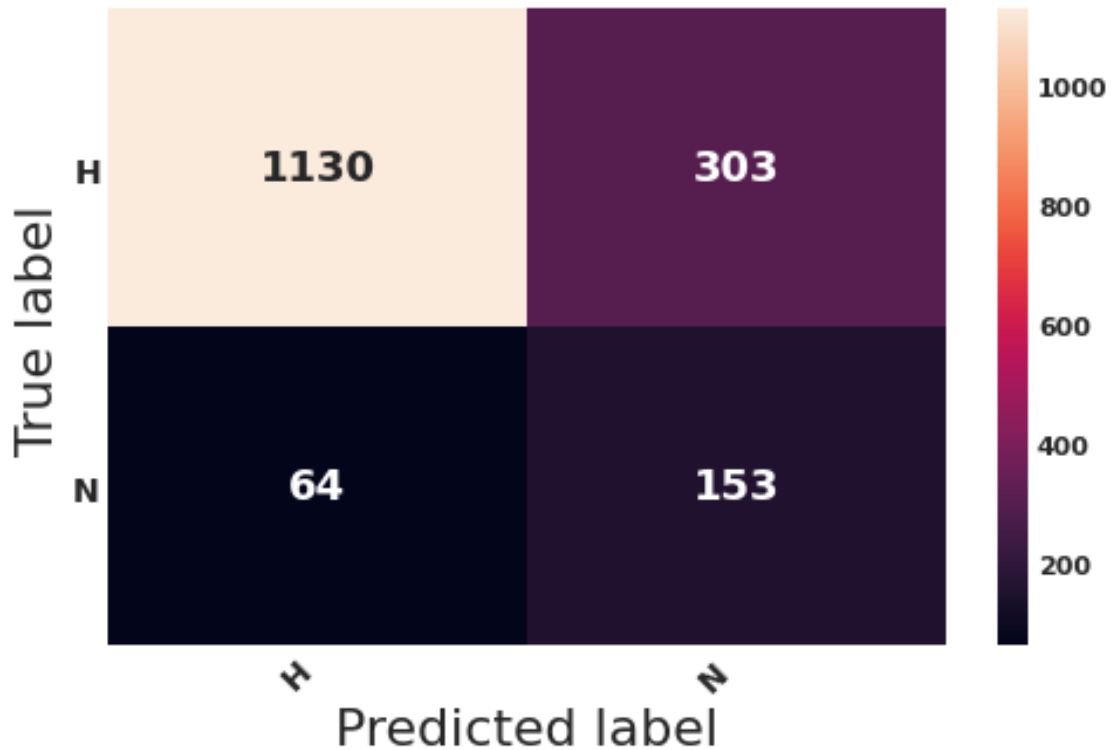
```
[91]: kNeighborsClassifier = KNeighborsClassifier(n_neighbors=1)
kNeighborsClassifier.fit(trainX,trainY)
KNCScores = cross_val_score(kNeighborsClassifier,trainX ,trainY,
    ↪cv=10,scoring='accuracy')
KNCScore = kNeighborsClassifier.score(testX,testY)
print("Accuracy Score : %f"%(KNCScore))
print("Accuracy Average Score of K Fold : %f"%(KNCScores.mean()))

knc_perdiction = kNeighborsClassifier.predict(testX)

knc_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(knc_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('KNeighborsClassifier',knc_con_matrix,train['label'].
    ↪unique())
```

Accuracy Score : 0.777576

Accuracy Average Score of K Fold : 0.760299



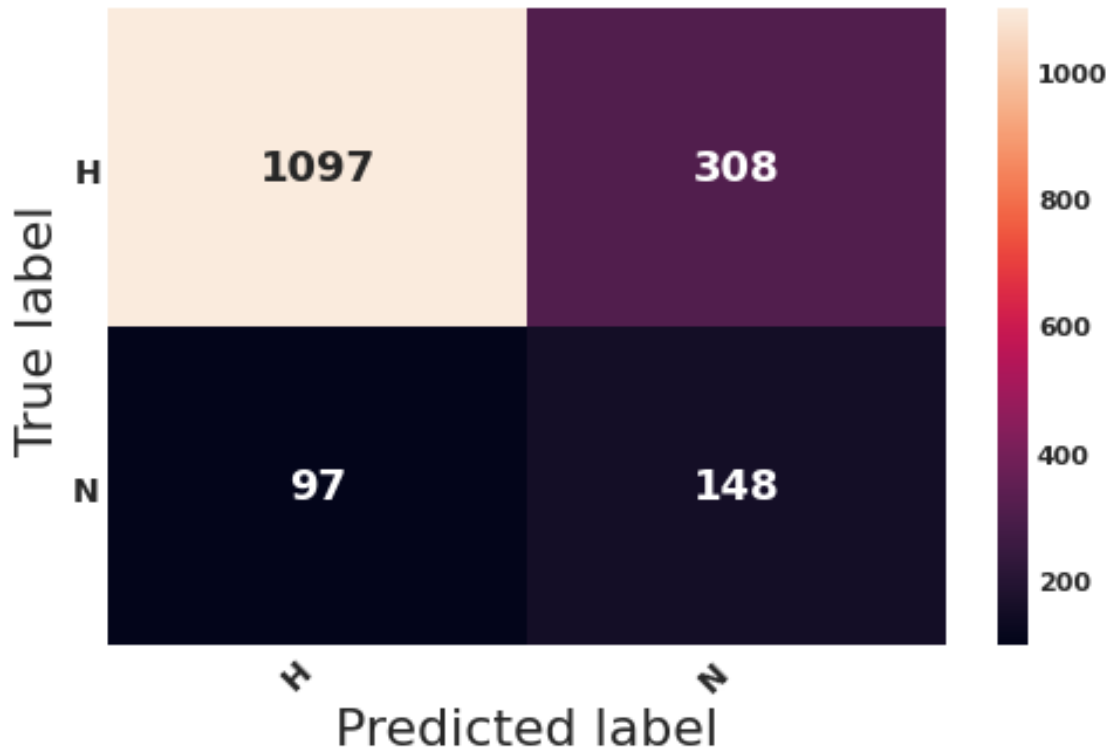
```
[92]: lda = LinearDiscriminantAnalysis()
lda.fit(trainX,trainY)
ldascores = cross_val_score(lda,trainX ,trainY, cv=10,scoring='accuracy')
ldascore = lda.score(testX,testY)
print("Accuracy Score : %f"%(ldascore))
print("Accuracy Average Score of K Fold : %f"%(ldascores.mean()))

lda_perdiction = lda.predict(testX)

lda_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(lda_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('LinearDiscriminantAnalysis',lda_con_matrix,train['label'].
    ↳unique())
```

Accuracy Score : 0.754545

Accuracy Average Score of K Fold : 0.732239



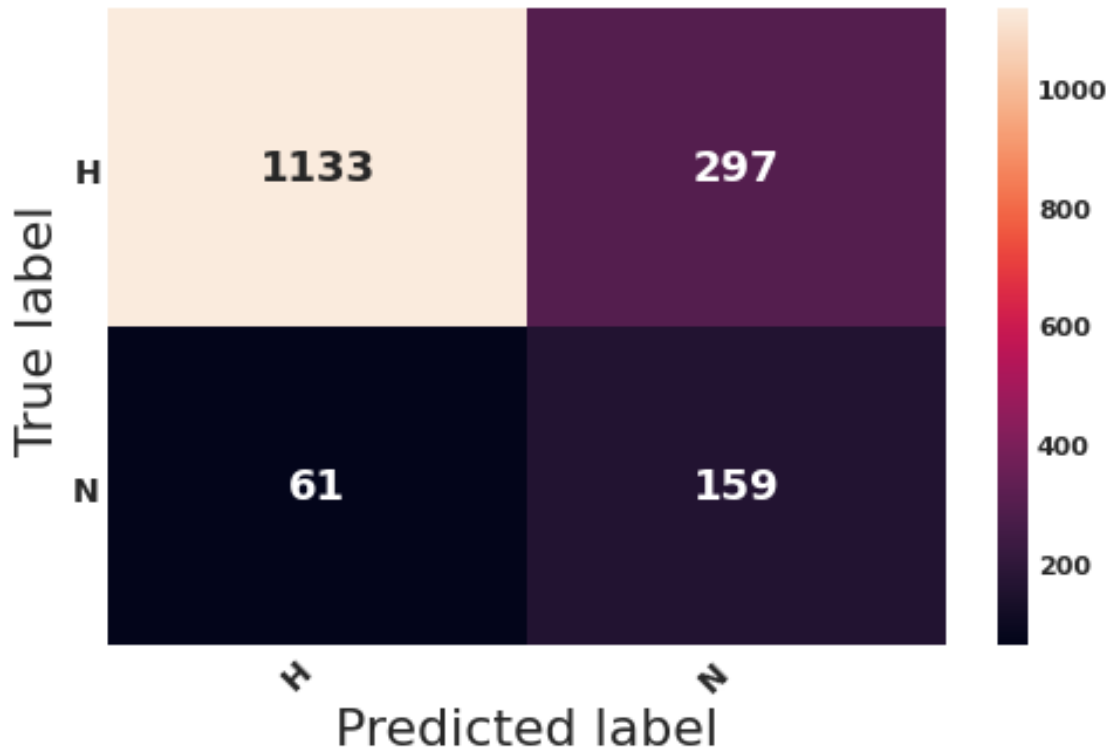
```
[93]: randomForestClassifier = RandomForestClassifier()
randomForestClassifier.fit(trainX,trainY)
RFCScores = cross_val_score(randomForestClassifier,trainX ,trainY,
    ↪cv=10,scoring='accuracy')
RFCScore = randomForestClassifier.score(testX,testY)
print("Accuracy Score : %f"%(RFCScore))
print("Accuracy Average Score of K Fold : %f"%(RFCScores.mean()))

rfc_perdiction = randomForestClassifier.predict(testX)

rfc_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(rfc_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('RandomForestClassifier',rfc_con_matrix,train['label'].
    ↪unique())
```

Accuracy Score : 0.783030

Accuracy Average Score of K Fold : 0.758209



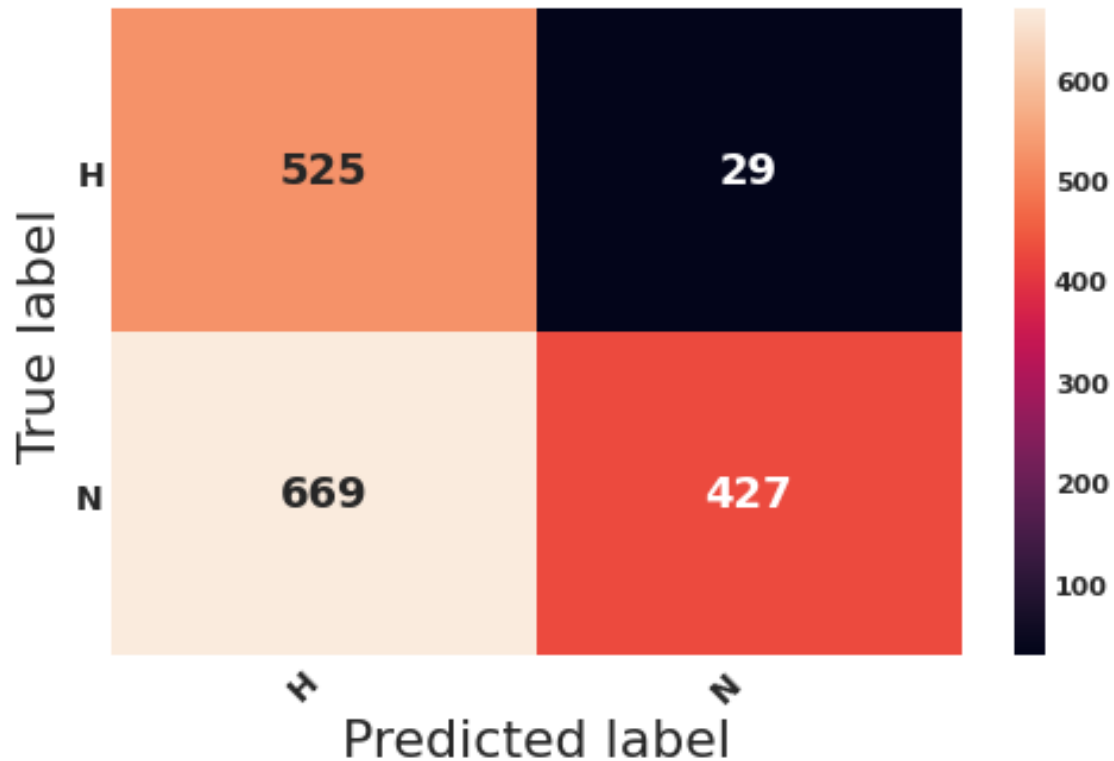
```
[94]: gaussianNB = GaussianNB()
gaussianNB.fit(trainX,trainY)
GNBScores = cross_val_score(gaussianNB,trainX ,trainY, cv=10,scoring='accuracy')
GNBScore = gaussianNB.score(testX,testY)
print("Accuracy Score : %f"%(GNBScore))
print("Accuracy Average Score of K Fold : %f"%(GNBScores.mean()))

GNB_perdiction = gaussianNB.predict(testX)

GNB_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(GNB_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('gaussianNB',GNB_con_matrix,train['label'].unique())
```

Accuracy Score : 0.576970

Accuracy Average Score of K Fold : 0.574030



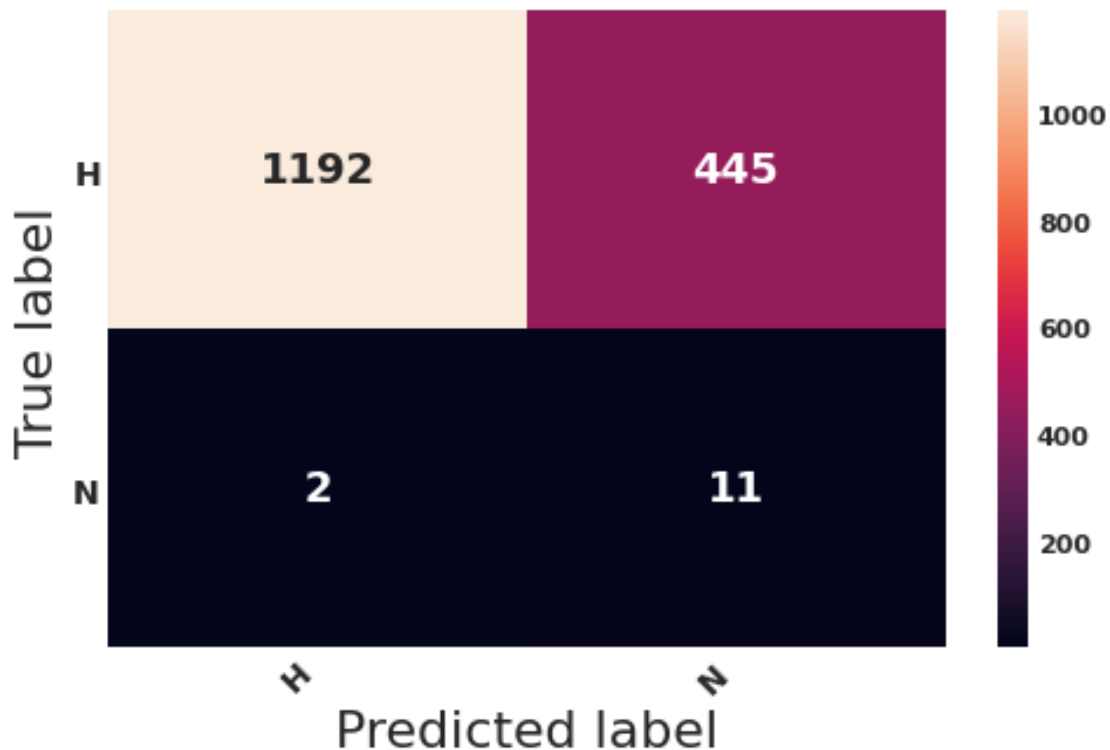
```
[95]: GBC = GradientBoostingClassifier()
GBC.fit(trainX,trainY)
GBCScores = cross_val_score(GBC,trainX ,trainY, cv=10,scoring='accuracy')
GBCScore = GBC.score(testX,testY)
print("Accuracy Score : %f"%(GBCScore))
print("Accuracy Average Score of K Fold : %f"%(GBCScores.mean()))

GBC_perdiction = GBC.predict(testX)

GBC_con_matrix = confusion_matrix(
    labelEncoder.inverse_transform(GBC_perdiction),
    labelEncoder.inverse_transform(testY),
    train['label'].unique()
)
print_confusion_matrix('GradientBoostingClassifier',GBC_con_matrix,train['label'].
    ↳unique())
```

Accuracy Score : 0.729091

Accuracy Average Score of K Fold : 0.725672



```
[96]: predictionModals = {
    'Accuracy': [LRScore, RFCScore, GBCTScore, BNBScore, DTCScore, KNCSTcore, ldscore, GNBSTcore, LSVCTScore],
    'Model': ['LR', 'RFC', 'GBC', 'BerNB', 'DTC', 'KNC', 'LinDA', 'GausNB', 'LinSVC']
}
predictionModelDF = pd.DataFrame(predictionModals, columns=['Accuracy', 'Model'])
print("detail performance of all parameters : ")
predictionModelDF
```

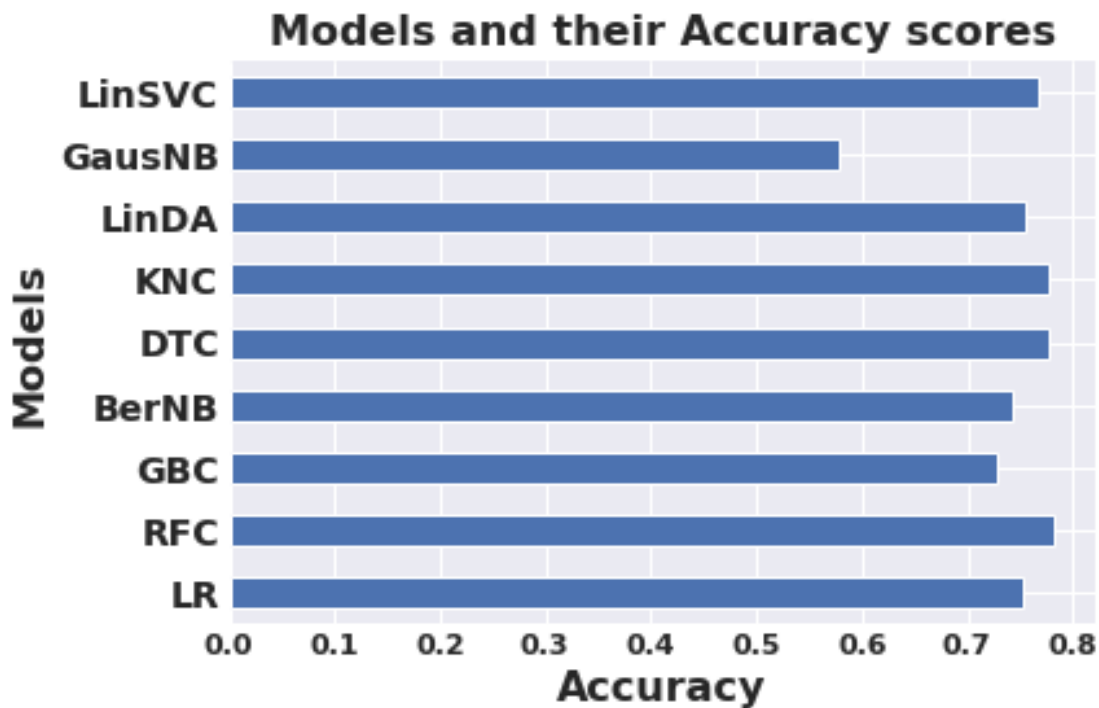
detail performance of all parameters :

```
[96]:
```

	Accuracy	Model
0	0.752121	LR
1	0.783030	RFC
2	0.729091	GBC
3	0.742424	BerNB
4	0.777576	DTC
5	0.777576	KNC
6	0.754545	LinDA
7	0.576970	GausNB
8	0.767879	LinSVC

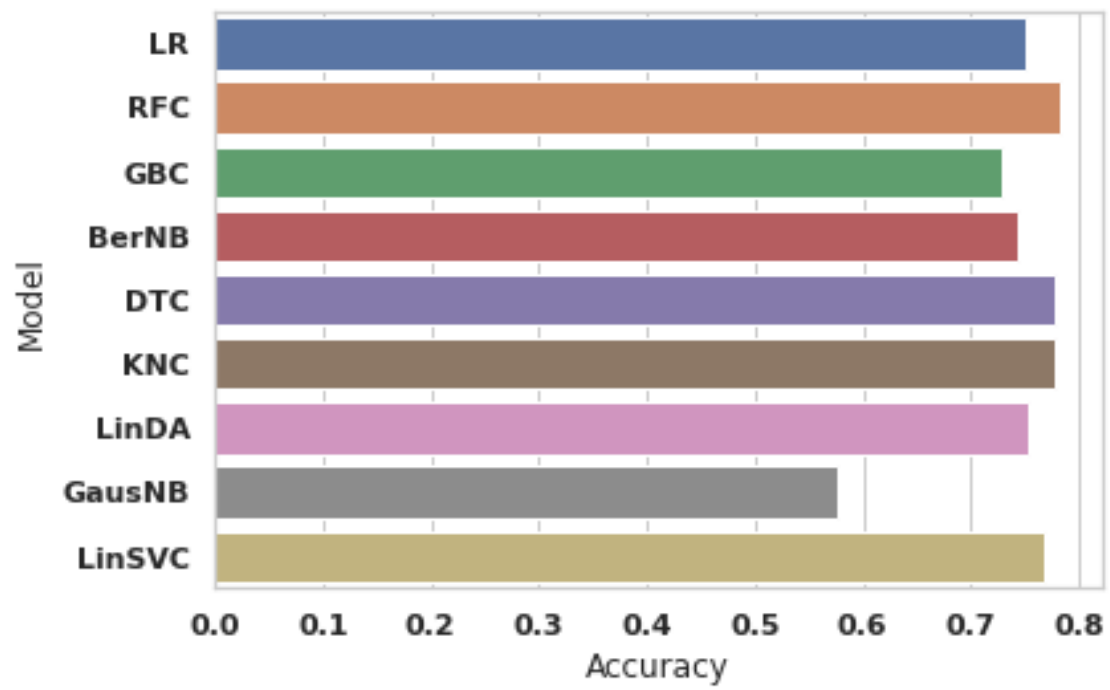


```
[97]: plt.rc('ytick', labels=14)
Plot = predictionModelDF.plot.barh(x='Model', y='Accuracy', legend=None)
Plot.set_ylabel('Models',fontdict={'fontsize': 16, 'fontweight': 'heavy'})
Plot.set_xlabel('Accuracy',fontdict={'fontsize': 16, 'fontweight': 'heavy'})
Plot.set_title("Models and their Accuracy scores",fontdict={'fontsize': 16,
↳'fontweight': 'heavy'})
fig = Plot.get_figure()
fig.savefig("ModelsAccuracy.
↳jpeg",dpi=600,format='jpeg',quality=100,pad_inches=0.1,
transparent=True, bbox_inches='tight')
```



```
[ ]:
```

```
[98]: sns.set_theme(style="whitegrid")
ax = sns.barplot(x="Accuracy", y="Model", data=predictionModelDF)
fig = ax.get_figure()
fig.savefig("ModelsAccuracy_seaborn.
↳jpeg",dpi=600,format='jpeg',quality=100,pad_inches=0.1,
transparent=True, bbox_inches='tight')
```



[ ]: