

sms-classifier

March 18, 2024

importing libraries

```
[8]: #Importing all the libraries to be used
import warnings
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pandas as pd
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from nltk.stem import WordNetLemmatizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.pipeline import Pipeline
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.model_selection import cross_val_score
from matplotlib.colors import ListedColormap
from sklearn.metrics import precision_score, recall_score, \
    plot_confusion_matrix, classification_report, accuracy_score, f1_score
from sklearn import metrics
```

loading data

```
[15]: data = pd.read_csv('spam.csv', encoding='latin-1')
```

```
[16]: data.head()
```

```
[16]:      v1      v2 Unnamed: 2 \
0  ham  Go until jurong point, crazy.. Available only ...  NaN
1  ham                Ok lar... Joking wif u oni...  NaN
2  spam  Free entry in 2 a wkly comp to win FA Cup fina...  NaN
3  ham  U dun say so early hor... U c already then say...  NaN
```

```
4    ham    Nah I don't think he goes to usf, he lives aro...    NaN
```

```
    Unnamed: 3 Unnamed: 4
0         NaN         NaN
1         NaN         NaN
2         NaN         NaN
3         NaN         NaN
4         NaN         NaN
```

```
[17]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0    v1              5572 non-null   object
1    v2              5572 non-null   object
2    Unnamed: 2      50 non-null     object
3    Unnamed: 3      12 non-null     object
4    Unnamed: 4       6 non-null     object
dtypes: object(5)
memory usage: 217.8+ KB
```

```
[18]: # Dropping the redundant looking collumns (for this project)
to_drop = ["Unnamed: 2","Unnamed: 3","Unnamed: 4"]
data = data.drop(data[to_drop], axis=1)
# Renaming the columns because I feel fancy today
data.rename(columns = {"v1":"Target", "v2":"Text"}, inplace = True)
data.head()
```

```
[18]:    Target    Text
0    ham    Go until jurong point, crazy.. Available only ...
1    ham                                Ok lar... Joking wif u oni...
2    spam    Free entry in 2 a wkly comp to win FA Cup fina...
3    ham    U dun say so early hor... U c already then say...
4    ham    Nah I don't think he goes to usf, he lives aro...
```

data expolaration

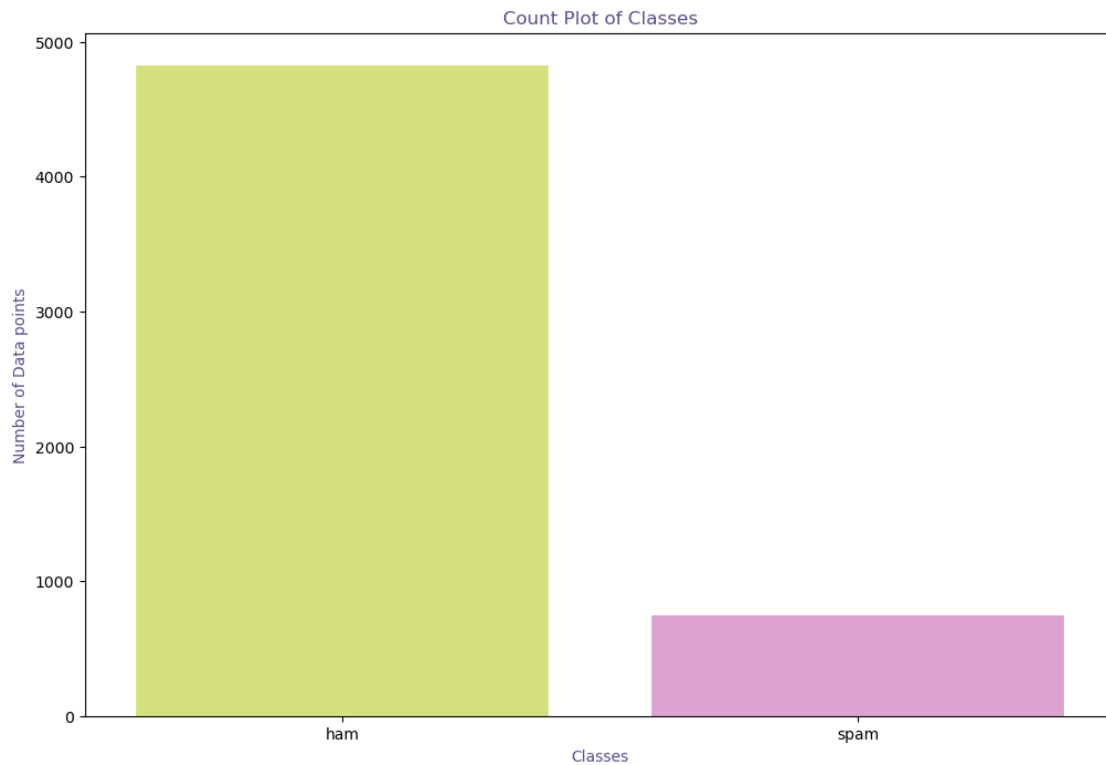
```
[19]: #Palette
cols= ["#E1F16B", "#E598D8"]
#first of all let us evaluate the target and find out if our data is imbalanced
↳ or not
plt.figure(figsize=(12,8))
fg = sns.countplot(x= data["Target"], palette= cols)
fg.set_title("Count Plot of Classes", color="#58508d")
```

```

fig.set_xlabel("Classes", color="#58508d")
fig.set_ylabel("Number of Data points", color="#58508d")

```

[19]: Text(0, 0.5, 'Number of Data points')



feature engineering

```

[21]: import nltk
      nltk.download('punkt')

```

```

[nltk_data] Downloading package punkt to C:\Users\CHANDRA
[nltk_data] ADITYA\AppData\Roaming\nltk_data...
[nltk_data] Unzipping tokenizers\punkt.zip.

```

[21]: True

```

[22]: #Adding a column of numbers of characters, words and sentences in each msg
data["No_of_Characters"] = data["Text"].apply(len)
data["No_of_Words"] = data.apply(lambda row: nltk.word_tokenize(row["Text"]),
                                ↪axis=1).apply(len)
data["No_of_sentence"] = data.apply(lambda row: nltk.sent_tokenize(row["Text"]),
                                ↪axis=1).apply(len)

```

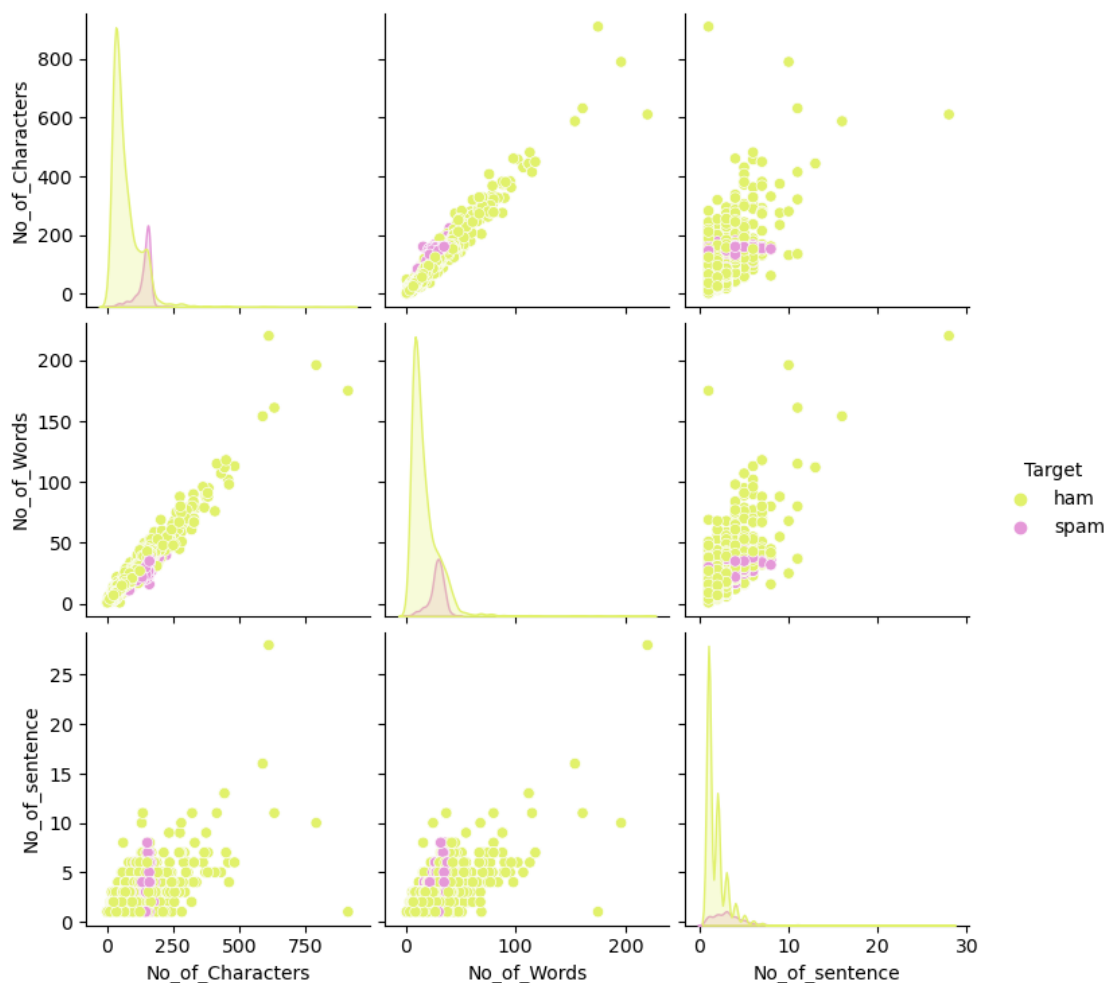
```
data.describe().T
```

```
[22]:
```

	count	mean	std	min	25%	50%	75%	max
No_of_Characters	5572.0	80.118808	59.690841	2.0	36.0	61.0	121.0	910.0
No_of_Words	5572.0	18.695621	13.742587	1.0	9.0	15.0	27.0	220.0
No_of_sentence	5572.0	1.970747	1.417778	1.0	1.0	1.0	2.0	28.0

```
[23]: plt.figure(figsize=(12,8))  
fg = sns.pairplot(data=data, hue="Target",palette=cols)  
plt.show(fg)
```

<Figure size 1200x800 with 0 Axes>



outlier detection

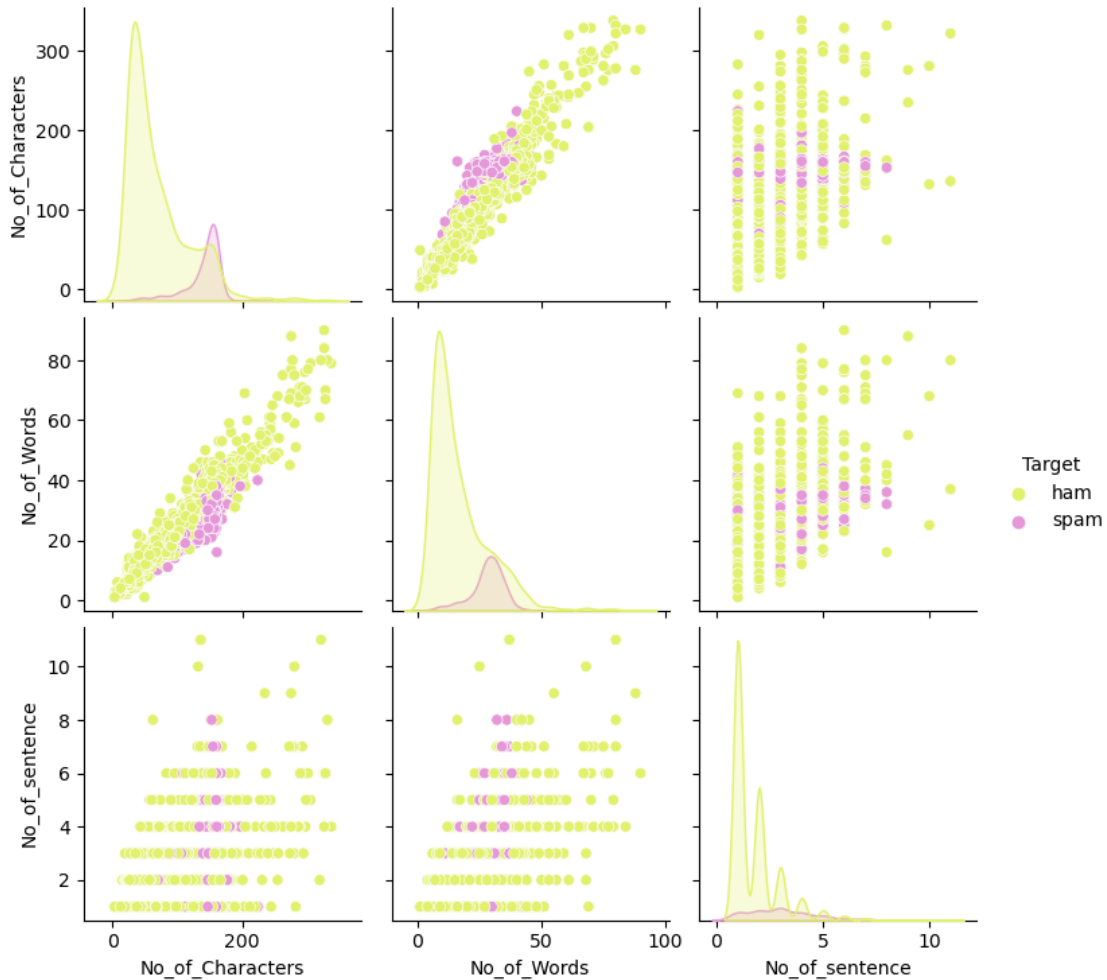
```
[24]: #Dropping the outliers.  
data = data[(data["No_of_Characters"]<350)]
```

```
data.shape
```

```
[24]: (5548, 5)
```

```
[25]: plt.figure(figsize=(12,8))  
fg = sns.pairplot(data=data, hue="Target",palette=cols)  
plt.show(fg)
```

<Figure size 1200x800 with 0 Axes>



data preprocessing

```
[26]: #Lets have a look at a sample of texts before cleaning  
print("\033[1m\u001b[45;1m The First 5 Texts:\033[0m",*data["Text"][:5], sep =_  
      ↪ "\n")
```

The First 5 Texts:

Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat...
 Ok lar... Joking wif u oni...
 Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry question(std txt rate)T&C's apply 08452810075over18's
 U dun say so early hor... U c already then say...
 Nah I don't think he goes to usf, he lives around here though

```
[27]: # Defining a function to clean up the text
def Clean(Text):
    sms = re.sub('[^a-zA-Z]', ' ', Text) #Replacing all non-alphabetic
    ↪characters with a space
    sms = sms.lower() #converting to lowercase
    sms = sms.split()
    sms = ' '.join(sms)
    return sms

data["Clean_Text"] = data["Text"].apply(Clean)
#Lets have a look at a sample of texts after cleaning
print("\033[1m\u001b[45;1m The First 5 Texts after cleaning:
    ↪\033[0m",*data["Clean_Text"][:5], sep = "\n")
```

The First 5 Texts after cleaning:

go until jurong point crazy available only in bugis n great world la e buffet
 cine there got amore wat
 ok lar joking wif u oni
 free entry in a wkly comp to win fa cup final tkts st may text fa to to receive
 entry question std txt rate t c s apply over s
 u dun say so early hor u c already then say
 nah i don t think he goes to usf he lives around here though

C:\Users\CHANDRA ADITYA\AppData\Local\Temp\ipykernel_2408\4237260299.py:9:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
data["Clean_Text"] = data["Text"].apply(Clean)
```

tokenization

```
[28]: data["Tokenize_Text"]=data.apply(lambda row: nltk.
    ↪word_tokenize(row["Clean_Text"]), axis=1)

print("\033[1m\u001b[45;1m The First 5 Texts after Tokenizing:
    ↪\033[0m",*data["Tokenize_Text"][:5], sep = "\n")
```

The First 5 Texts after Tokenizing:

```
['go', 'until', 'jurong', 'point', 'crazy', 'available', 'only', 'in', 'bugis',
'n', 'great', 'world', 'la', 'e', 'buffet', 'cine', 'there', 'got', 'amore',
'wat']
['ok', 'lar', 'joking', 'wif', 'u', 'oni']
['free', 'entry', 'in', 'a', 'wkly', 'comp', 'to', 'win', 'fa', 'cup', 'final',
'tkts', 'st', 'may', 'text', 'fa', 'to', 'to', 'receive', 'entry', 'question',
'std', 'txt', 'rate', 't', 'c', 's', 'apply', 'over', 's']
['u', 'dun', 'say', 'so', 'early', 'hor', 'u', 'c', 'already', 'then', 'say']
['nah', 'i', 'don', 't', 'think', 'he', 'goes', 'to', 'usf', 'he', 'lives',
'around', 'here', 'though']
```

C:\Users\CHANDRA ADITYA\AppData\Local\Temp\ipykernel_2408\3712914543.py:1:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
data["Tokenize_Text"]=data.apply(lambda row:
nltk.word_tokenize(row["Clean_Text"]), axis=1)
```

removing stopwords

```
[30]: import nltk
nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to C:\Users\CHANDRA
[nltk_data] ADITYA\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

[30]: True

```
[31]: # Removing the stopwords function
def remove_stopwords(text):
    stop_words = set(stopwords.words("english"))
    filtered_text = [word for word in text if word not in stop_words]
    return filtered_text

data["Nostopword_Text"] = data["Tokenize_Text"].apply(remove_stopwords)

print("\033[1m\u001b[45;1m The First 5 Texts after removing the stopwords:
↪\033[0m",*data["Nostopword_Text"][:5], sep = "\n")
```

The First 5 Texts after removing the stopwords:

```
['go', 'jurong', 'point', 'crazy', 'available', 'bugis', 'n', 'great', 'world',
'la', 'e', 'buffet', 'cine', 'got', 'amore', 'wat']
['ok', 'lar', 'joking', 'wif', 'u', 'oni']
['free', 'entry', 'wkly', 'comp', 'win', 'fa', 'cup', 'final', 'tkts', 'st',
'may', 'text', 'fa', 'receive', 'entry', 'question', 'std', 'txt', 'rate', 'c',
```

```
'apply']
['u', 'dun', 'say', 'early', 'hor', 'u', 'c', 'already', 'say']
['nah', 'think', 'goes', 'usf', 'lives', 'around', 'though']
```

C:\Users\CHANDRA ADITYA\AppData\Local\Temp\ipykernel_2408\3551966202.py:7:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
data["Nostopword_Text"] = data["Tokenize_Text"].apply(remove_stopwords)
```

lemmatization

```
[35]: import nltk
nltk.download('wordnet')
nltk.download('omw-1.4')
```

[nltk_data] Downloading package wordnet to C:\Users\CHANDRA

[nltk_data] ADITYA\AppData\Roaming\nltk_data...

[nltk_data] Package wordnet is already up-to-date!

[nltk_data] Downloading package omw-1.4 to C:\Users\CHANDRA

[nltk_data] ADITYA\AppData\Roaming\nltk_data...

[35]: True

```
[36]: lemmatizer = WordNetLemmatizer()
# lemmatize string
def lemmatize_word(text):
    #word_tokens = word_tokenize(text)
    # provide context i.e. part-of-speech
    lemmas = [lemmatizer.lemmatize(word, pos='v') for word in text]
    return lemmas

data["Lemmatized_Text"] = data["Nostopword_Text"].apply(lemmatize_word)
print("\033[1m\u001b[45;1m The First 5 Texts after lemitization:
↪\033[0m",*data["Lemmatized_Text"][:5], sep = "\n")
```

The First 5 Texts after lemitization:

```
['go', 'jurong', 'point', 'crazy', 'available', 'bugis', 'n', 'great', 'world',
'la', 'e', 'buffet', 'cine', 'get', 'amore', 'wat']
```

```
['ok', 'lar', 'joke', 'wif', 'u', 'oni']
```

```
['free', 'entry', 'wkly', 'comp', 'win', 'fa', 'cup', 'final', 'tkts', 'st',
'may', 'text', 'fa', 'receive', 'entry', 'question', 'std', 'txt', 'rate', 'c',
'apply']
```

```
['u', 'dun', 'say', 'early', 'hor', 'u', 'c', 'already', 'say']
```

```
['nah', 'think', 'go', 'usf', 'live', 'around', 'though']
```

vectorize


```
[37]: #Creating a corpus of text feature to encode further into vectorized form
corpus= []
for i in data["Lemmatized_Text"]:
    msg = ' '.join([row for row in i])
    corpus.append(msg)

corpus[:5]
print("\033[1m\u001b[45;1m The First 5 lines in corpus :\033[0m",*corpus[:5],
      ↪sep = "\n")
```

The First 5 lines in corpus :

```
go jurong point crazy available bugis n great world la e buffet cine get amore
wat
ok lar joke wif u oni
free entry wkly comp win fa cup final tkts st may text fa receive entry question
std txt rate c apply
u dun say early hor u c already say
nah think go usf live around though
```

```
[38]: #Changing text data in to numbers.
tfidf = TfidfVectorizer()
X = tfidf.fit_transform(corpus).toarray()
#Let's have a look at our feature
X.dtype
```

```
[38]: dtype('float64')
```

```
[39]: #Label encode the Target and use it as y
label_encoder = LabelEncoder()
data["Target"] = label_encoder.fit_transform(data["Target"])
```

model building

```
[40]: #Setting values for labels and feature as y and X(we already did X in
      ↪vectorizing...)
y = data["Target"]
# Splitting the testing and training sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
      ↪random_state=42)
```

```
[41]: #Testing on the following classifiers
classifiers = [MultinomialNB(),
                RandomForestClassifier(),
                KNeighborsClassifier(),
                SVC()]
for cls in classifiers:
    cls.fit(X_train, y_train)
```

```
# Dictionary of pipelines and model types for ease of reference
pipe_dict = {0: "NaiveBayes", 1: "RandomForest", 2: "KNeighbours", 3: "SVC"}
```

```
[42]: # Cossvalidation
for i, model in enumerate(classifiers):
    cv_score = cross_val_score(model, X_train, y_train, scoring="accuracy", cv=10)
    print("%s: %f " % (pipe_dict[i], cv_score.mean()))
```

```
NaiveBayes: 0.967552
RandomForest: 0.975437
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will

change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-

packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-

packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-

packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-

packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

C:\Users\CHANDRA ADITYA\anaconda3\lib\site-

packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

KNeighbours: 0.911450

SVC: 0.974086

evaluating models

```
[ ]: # Model Evaluation
# creating lists of varios scores
precision = []
recall = []
f1_score = []
trainset_accuracy = []
testset_accuracy = []

for i in classifiers:
    pred_train = i.predict(X_train)
    pred_test = i.predict(X_test)
    prec = metrics.precision_score(y_test, pred_test)
    recal = metrics.recall_score(y_test, pred_test)
    f1_s = metrics.f1_score(y_test, pred_test)
    train_accuracy = model.score(X_train,y_train)
    test_accuracy = model.score(X_test,y_test)

    #Appending scores
    precision.append(prec)
    recall.append(recal)
    f1_score.append(f1_s)
    trainset_accuracy.append(train_accuracy)
    testset_accuracy.append(test_accuracy)
```

```
[ ]: # initialise data of lists.
data = {'Precision':precision,
'Recall':recall,
'F1score':f1_score,
'Accuracy on Testset':testset_accuracy,
'Accuracy on Trainset':trainset_accuracy}
# Creates pandas DataFrame.
Results = pd.DataFrame(data, index=["NaiveBayes", "RandomForest",
↪ "KNeighbours", "SVC"])
```

```
[ ]: cmap2 = ListedColormap(["#E2CCFF", "#E598D8"])
Results.style.background_gradient(cmap=cmap2)
```

```
[ ]: cmap = ListedColormap(["#E1F16B", "#E598D8"])
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(15,10))

for cls, ax in zip(classifiers, axes.flatten()):
    plot_confusion_matrix(cls,
                           X_test,
                           y_test,
                           ax=ax,
                           cmap= cmap,
                           )
```

```
ax.title.set_text(type(cls).__name__)  
plt.tight_layout()  
plt.show()
```