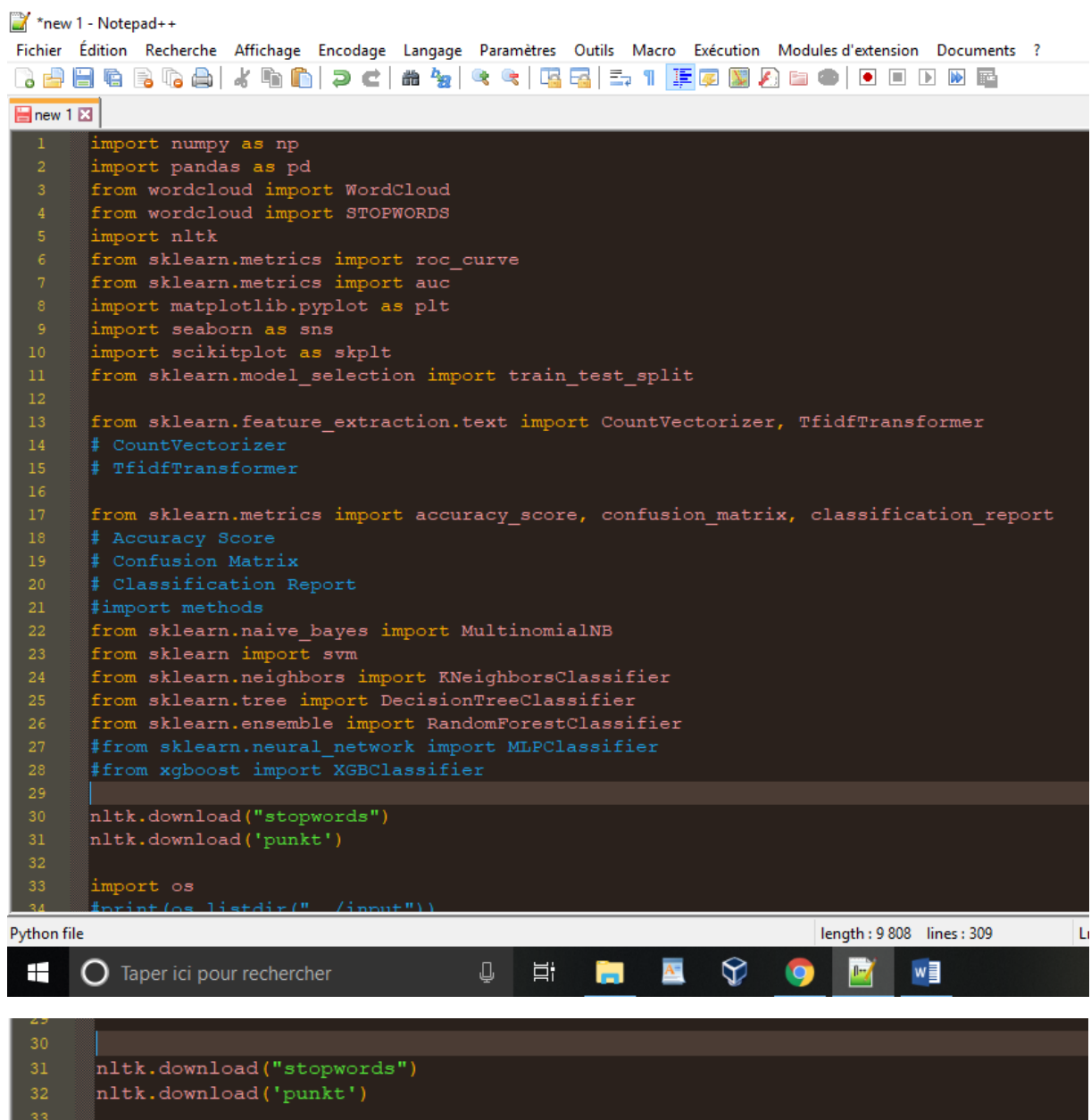


Spam Project of Machine Learning:

Members:

- **Edna Batana Diane**
- **Loukrimi Souad**
- **Zemmouri Kenza**



```
*new 1 - Notepad++
Fichier  Édition  Recherche  Affichage  Encodage  Langage  Paramètres  Outils  Macro  Exécution  Modules d'extension  Documents  ?

new 1 x
1  import numpy as np
2  import pandas as pd
3  from wordcloud import WordCloud
4  from wordcloud import STOPWORDS
5  import nltk
6  from sklearn.metrics import roc_curve
7  from sklearn.metrics import auc
8  import matplotlib.pyplot as plt
9  import seaborn as sns
10 import scikitplot as skplt
11 from sklearn.model_selection import train_test_split
12
13 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
14 # CountVectorizer
15 # TfidfTransformer
16
17 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
18 # Accuracy Score
19 # Confusion Matrix
20 # Classification Report
21 #import methods
22 from sklearn.naive_bayes import MultinomialNB
23 from sklearn import svm
24 from sklearn.neighbors import KNeighborsClassifier
25 from sklearn.tree import DecisionTreeClassifier
26 from sklearn.ensemble import RandomForestClassifier
27 #from sklearn.neural_network import MLPClassifier
28 #from xgboost import XGBClassifier
29
30 nltk.download("stopwords")
31 nltk.download('punkt')
32
33 import os
34 #print(os.listdir("/input"))

Python file  length: 9 808  lines: 309  L
Taper ici pour rechercher
```

```

33
34 import os
35 #print(os.listdir("../input"))
36

```

Read Data :

```

39
40 # ##### Read Data:
41 df = pd.read_csv("/Users/elizabethlorelei/Downloads/spam.csv", encoding = 'latin-1')
42

```

Show Data :

```

43
44
45 # ##### Show Data :
46 df.head(11)
47

```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only ...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
5	spam	FreeMsg Hey there darling it's been 3 week's n...	NaN	NaN	NaN
6	ham	Even my brother is not like to speak with me. ...	NaN	NaN	NaN
7	ham	As per your request 'Melle Melle (Oru Minnamin...	NaN	NaN	NaN
8	spam	WINNER!! As a valued network customer you have...	NaN	NaN	NaN
9	spam	Had your mobile 11 months or more? U R entitle...	NaN	NaN	NaN
10	ham	I'm gonna be home soon and i don't want to tal...	NaN	NaN	NaN

```

63
64 df = df.drop(columns=['Unnamed: 2','Unnamed: 3','Unnamed: 4'])
65 df.columns = ['Label', 'Message']
66 df.head()
67

```

	Label	Message
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...

Splitting the labels and the data separately:

```

77
78 # ##### Splitting the labels and the data separately :
79 df_labels = df['Label']
80 df_labels.head(11)
81

```

```

0      ham
1      ham
2     spam
3      ham
4      ham
5     spam
6      ham
7      ham
8     spam
9     spam
10     ham
Name: Label, dtype: object

```

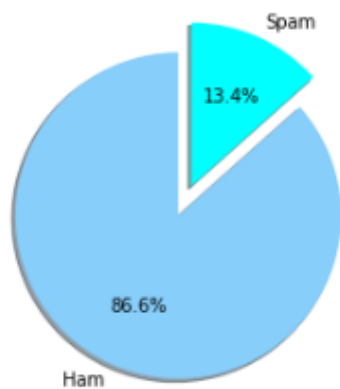
Data Visualization:

- To check the most used word in Ham sms and Spam SMS
- To visualize the percentage of Ham and Spam SMS

[illegible][illegible]

Plotting ham and spam data % in pie chart:

```
145 # ##### Plotting ham and spam data % in pie chart :
146
147 count_Class = pd.value_counts(df.Label, sort = True)
148
149 # Data to Plot
150 labels = 'Ham', 'Spam'
151 sizes = [count_Class[0], count_Class[1]]
152 colors = ['lightskyblue', 'aqua']
153 explode = (0.1, 0.1)
154
155 # Plot
156 plt.pie(sizes, explode = explode, labels = labels, colors = colors,
157         autopct = '%1.1f%%', shadow = True, startangle = 90)
158 plt.axis('equal')
159 plt.show()
```



Splitting the Test and Train Data:

```
165
166 # ##### Splitting the Test and Train Data:
167
168 train_set, test_set, train_label, test_label = train_test_split(df, df_labels, test_size = 0.33, random_state = 42)
169 print(train_set.shape)
170 print(test_set.shape)
171 print("\nThe Trainset consists of {} records and {} features".format(train_set.shape[0],train_set.shape[1]))
172 print("\nThe Testset consists of {} records and {} features".format(test_set.shape[0],train_set.shape[1]))
173
174
175
```

(3733, 2)

(1839, 2)

The Trainset consists of 3733 records and 2 features

The Testset consists of 1839 records and 2 features

Extracting N-grams from the Text Data:

```
176
177 # ##### Extracting N-grams from the Text Data:
178
179 countvect = CountVectorizer(ngram_range = (2,2), )
180 x_counts = countvect.fit(train_set.Message)
181
182 # preparing for training set
183 x_train_df = countvect.transform(train_set.Message)
184
185 # preparing for test set
186 x_test_df = countvect.transform(test_set.Message)
187
188
```

Data Model:

The Algorithms used below are:

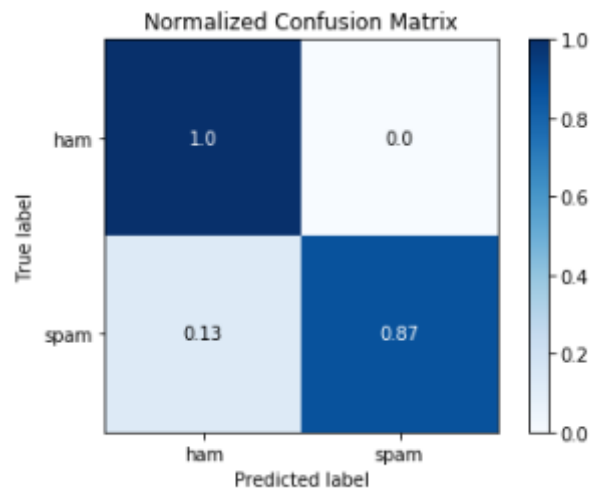
- Naive Bayes
- K-Nearest
- Decision Tree
- Support Vector Machine
- Random Forest

Naive Bayes classifier :

```
190
191 # ##### Naive Bayes classifier :
192
193
194 clf = MultinomialNB()
195 clf.fit(x_train_df, train_set.Label)
196 predicted_values_NB = clf.predict(x_test_df)
197 predictions = dict()
198 accuracy = accuracy_score(test_set.Label, predicted_values_NB)
199 predictions['Naive Bayes'] = accuracy * 100
200 confusionmatrix = confusion_matrix(test_set.Label, predicted_values_NB)
201 print("The accuracy of Naive Bayes classifier is {}".format(accuracy * 100))
202 print("\n", confusionmatrix)
203 skplt.metrics.plot_confusion_matrix(test_set.Label, predicted_values_NB, normalize = True)
204 plt.show()
205
```

The accuracy of Naive Bayes clasifier is 97.87928221859707%

```
[[1581    6]
 [33    219]]
```



K-Nearest Neighbors algorithm :

```

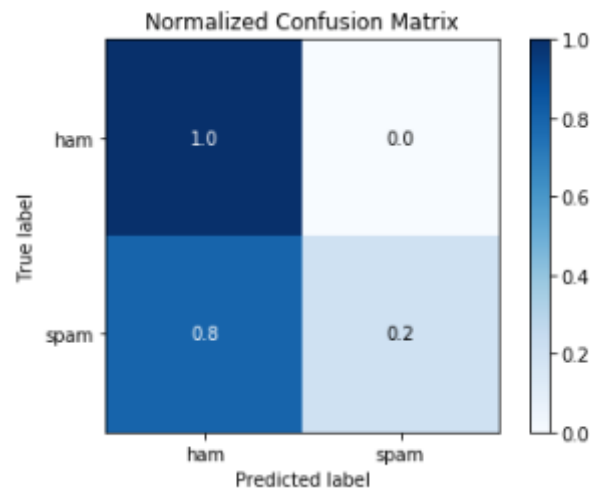
209 #KNN = KNeighborsClassifier(metric = 'euclidean')
210 KNN = KNeighborsClassifier()
211 KNN.fit(x_train_df, train_set.Label)
212 predicted_values_KNN = KNN.predict(x_test_df)
213 print(predicted_values_KNN)
214 accuracy_KNN = accuracy_score(test_set.Label, predicted_values_KNN)
215 predictions['K-Nearest Neighbors algorithm'] = accuracy_KNN * 100
216 print("\nThe accuracy of K-Nearest Neighbors algorithm is {}".format(accuracy_KNN * 100))
217 confusion_matrix_KNN = confusion_matrix(test_set.Label, predicted_values_KNN)
218 print("\n", confusion_matrix_KNN)
219 skplt.metrics.plot_confusion_matrix(test_set.Label, predicted_values_KNN, normalize = True)
220 plt.show()
221

```

```
['ham' 'ham' 'ham' ... 'ham' 'ham' 'ham']
```

The accuracy of K-Nearest Neighbors algorithm is 89.07014681892332%

```
[[1587    0]
 [201    51]]
```

Decision Tree learning :

```

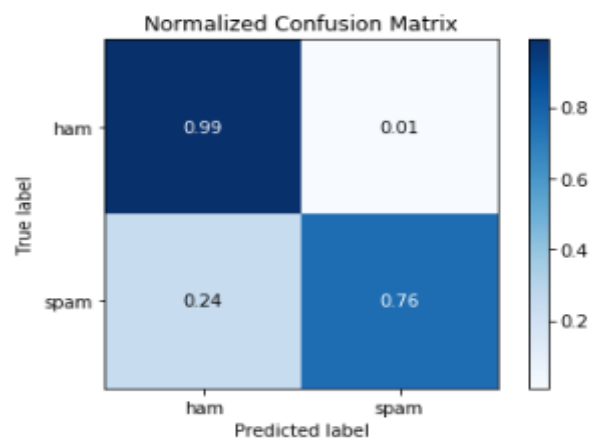
221
222 # ### Decision Tree learning :
223
224 DT = DecisionTreeClassifier()
225 DT.fit(x_train_df, train_set.Label)
226 predicted_values_DT = DT.predict(x_test_df)
227 print(predicted_values_DT)
228 accuracy_DT = accuracy_score(test_set.Label, predicted_values_DT)
229 predictions['Decision Tree learning'] = accuracy_DT * 100
230 print("\nThe accuracy of Decision Tree learning is {}".format(accuracy_DT * 100))
231 confusion_matrix_DT = confusion_matrix(test_set.Label, predicted_values_DT)
232 print("\n", confusion_matrix_DT)
233 skplt.metrics.plot_confusion_matrix(test_set.Label, predicted_values_DT, normalize = True)
234 plt.show()
235
236

```

['ham' 'ham' 'spam' ... 'ham' 'ham' 'spam'] The accuracy of Decision Tree learning is 96.08482871125612%

[[1576 11]

[61 191]]



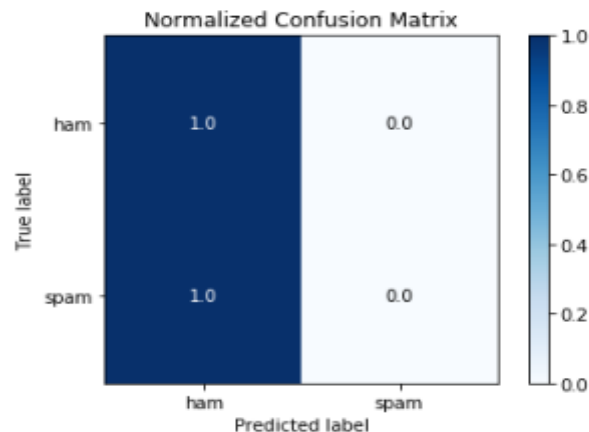
Support Vector Machine (SVM) :

```
238
239 # ##### Support Vector Machine (SVM) :
240
241 SVM = svm.SVC()
242 SVM.fit(x_train_df, train_set.Label)
243 predicted_values_SVM = SVM.predict(x_test_df)
244 print(predicted_values_SVM)
245 accuracy_SVM = accuracy_score(test_set.Label, predicted_values_SVM)
246 predictions['Support Vector Machine (SVM)'] = accuracy_SVM * 100
247 print("\nThe accuracy of Support Vector Machine (SVM) is {}".format(accuracy_SVM * 100))
248 confusion_matrix_SVM = confusion_matrix(test_set.Label, predicted_values_SVM)
249 print("\n", confusion_matrix_SVM)
250 skplt.metrics.plot_confusion_matrix(test_set.Label, predicted_values_SVM, normalize = True)
251 plt.show()
252
```

['ham' 'ham' 'ham' ... 'ham' 'ham' 'ham']

The accuracy of Support Vector Machine (SVM) is 86.2969004893964%

```
[[1587    0]
 [252     0]]
```

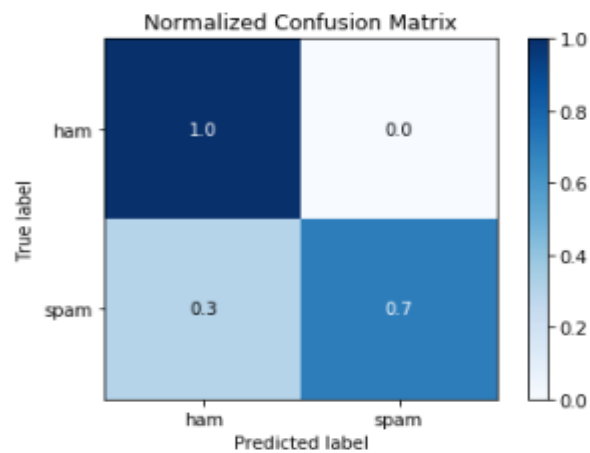


Random Forest:

```
253
254
255 # ##### Random Forest :
256
257 RF = RandomForestClassifier(n_estimators = 100, oob_score = True, random_state = 123456)
258 # n_estimators - количество деревьев в лесу
259 # oob_score - использовать ли образцы вне примеров для оценки точности обобщения
260 RF.fit(x_train_df, train_set.Label)
261 predicted_values_RF = RF.predict(x_test_df)
262 print(predicted_values_RF)
263 accuracy_RF = accuracy_score(test_set.Label, predicted_values_RF)
264 predictions['Random Forest'] = accuracy_RF * 100
265 print("\nThe accuracy of Random Forest is {}".format(accuracy_RF * 100))
266 confusion_matrix_RF = confusion_matrix(test_set.Label, predicted_values_RF)
267 print("\n", confusion_matrix_RF)
268 skplt.metrics.plot_confusion_matrix(test_set.Label, predicted_values_RF, normalize = True)
269 plt.show()
270
```

['ham' 'ham' 'ham' ... 'ham' 'ham' 'ham'] The accuracy of Random Forest is 95.86731919521479%

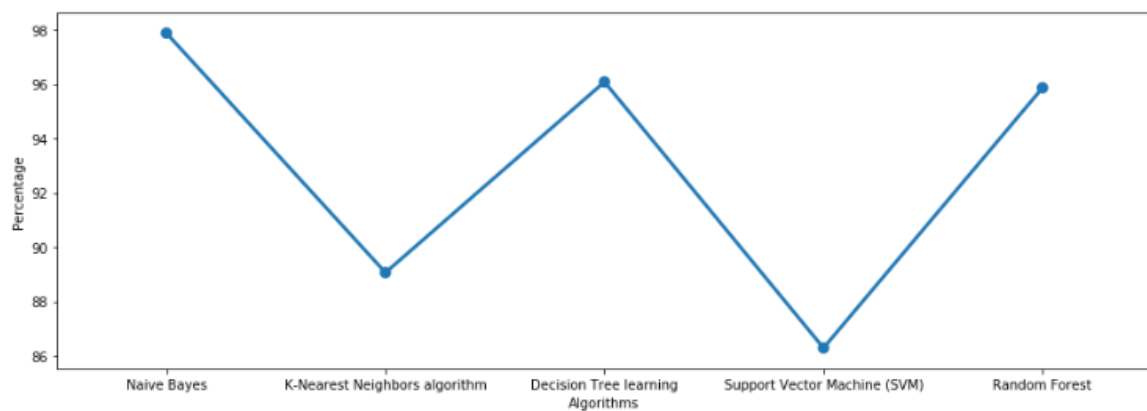
```
[[1587 0]
 [76  176]]
```



Méthodes Comparison :

```
271
272 # ##### Method Comparison:
273
274 fig, (ax1) = plt.subplots(ncols = 1, sharey = True, figsize = (15,5))
275 df = pd.DataFrame(list(predictions.items()), columns = ['Algorithms', 'Percentage'])
276 display(df)
277 sns.pointplot(x = "Algorithms", y = "Percentage", data = df, ax = ax1);
278
279
```

	Algorithms	Percentage
0	Naive Bayes	97.879282
1	K-Nearest Neighbors algorithm	89.070147
2	Decision Tree learning	96.084829
3	Support Vector Machine (SVM)	86.296900
4	Random Forest	95.867319



ROC Accuracy :

```

283 #pr, tpr, thresholds = roc_curve(testset.v1,predicted_values_XGB, pos_label=2)
284 test_prediction = test_set.Label.tolist()
285 predicted_values = predicted_values_NB.tolist()
286 test_prediction = [1 if pred=="spam" else 0 for pred in test_prediction]
287 predicted_values = [1 if pred=="spam" else 0 for pred in predicted_values]
288 fpr, tpr, thresholds = roc_curve(test_prediction,predicted_values)
289 roc_auc = auc(fpr, tpr)
290 print("The ROC Accuracy is {}".format(roc_auc))
291

```

The ROC Accuracy is 0.9326334503555676

```

296 plt.title('Receiver Operating Characteristic')
297 plt.plot(fpr, tpr, 'b',
298 label='AUC = %0.2f'% roc_auc)
299 plt.legend(loc='lower right')
300 plt.plot([0,1],[0,1],'r--')
301 plt.xlim([-0.1,1.2])
302 plt.ylim([-0.1,1.2])
303 plt.ylabel('True Positive Rate')
304 plt.xlabel('False Positive Rate')
305 plt.show()
306
307

```

