

Construction des classifieurs

Bachar RIMA

Emile YOUSSEF

Tasnim SHAQURA

April 22, 2019

Contents

1	Résumé du projet	2
2	Aperçu du module utility_ML.py	2
3	Importation des modules et configuration de l'environnement de travail	9
4	Importation des datasets	10
5	Pré-traitement des données	11
5.1	Pré-traitement d'un document	11
5.2	Pré-traitement des tokens	12
6	Visualisation des données	12
7	Vectorisation	14
8	Cross-validation scores	14
8.1	Résultats de la cross-validation	17
9	Calibrage des hyperparamètres des modèles en utilisant un GridSearch	17
9.1	Résultats du calibrage avec un GridSearch	19
10	Création d'un Pipeline pour le modèle Logistic Regression	20
11	Création d'un Pipeline pour le modèle Gaussian Naive Bayes	22

1 Résumé du projet

- Un jeu de données textuelles est mis à disposition sur Moodle.
- Il s'agit d'un corpus de 10.000 documents contenant des avis d'internautes sur des films.
- A chaque document est associé sa polarité selon l'avis (+1 : positif, -1 : négatif).
- Le fichier des documents est formaté dans un tableau csv (un avis par ligne), un autre fichier csv contient les polarités d'avis par document (-1/+1).
- Une correspondance directe existe entre les numéros des lignes des documents et des polarités.

2 Aperçu du module utility_ML.py

Le module `utility_ML.py` est utilisé pour externaliser les éléments de programmation utilitaires utilisées tout au long de notre projet. Ceci inclut :

1. Des **constantes** désignant des chemins vers les **ressources utilisées/créées** par le projet (*datasets, fichiers pickle, etc.*)
2. Des **constantes** désignant des **objets** utilisés dans le code du projet (*e.g. STOP_WORDS, POS_TAG_MAP, valeurs de paramétrage des fonctions utilisés dans le projet, etc.*)
3. Des **fonctions wrapper** pour les fonctions d'**importation**, de **fusion** et de **mélange** de datasets
4. Des **fonctions de prétraitement** (*e.g. remplacement des contractions, suppression de ponctuation, lemmatization, etc.*)
5. Des **classes utilitaires** utilisées dans le projet (*e.g. classe encapsulant les résultats d'une recherche par GridSearch*)

```
1 #MODULES
2 import re
3 import unicodedata
4 import warnings
5
6 import pandas as pd
7 import numpy as np
8 import contractions
9 import inflect
10
11 from collections import defaultdict
12
13 from nltk import pos_tag
14 from nltk import punkt
15 from nltk.corpus import stopwords, wordnet as wn
16 from nltk.tokenize import word_tokenize
```

```

17 from nltk.stem import WordNetLemmatizer
18
19 from sklearn.utils import shuffle
20 from sklearn.base import TransformerMixin
21
22 warnings.filterwarnings('ignore', category=FutureWarning)
23
24 #DEFINITION OF SOME OF THE RESOURCES USED THROUGHOUT THIS NOTEBOOK
25 DATA_PATH = 'Datasets/dataset.csv'
26 TARGET_PATH = 'Datasets/labels.csv'
27 TEST_DATA_PATH = 'Datasets/test_data.csv'
28 TEST_TARGET_PATH = 'Datasets/test_labels.csv'
29 IMDB_DATA_PATH = 'Datasets/imdb_reviews_scores.csv'
30 GAUSSIANNB_PATH = 'Classifiers/gaussianNB.pkl'
31 LOGISTICREGRESSION_PATH = 'Classifiers/logisticRegression.pkl'
32
33 #stop words set adapted to the context of the dataset
34 STOP_WORDS = set(stopwords.words('english'))
35 STOP_WORDS_EXCEPTIONS = set(('no', 'nor', 'not',))
36 STOP_WORDS_ADDITIONS = [
37     'film', 'films', 'filmed',
38     'movie', 'movies',
39     'character', 'characters',
40     'story', 'stories',
41     'scene', 'scenes',
42     'actor', 'actors', 'actress', 'actresses', 'act', 'acts', 'acted',
43     'acting'
44     'direct', 'directs', 'directed', 'directing', 'director', 'directors',
45     'script', 'scripts',
46     'plot', 'plots'
47 ]
48 STOP_WORDS.update(STOP_WORDS_ADDITIONS)
49 STOP_WORDS = STOP_WORDS - STOP_WORDS_EXCEPTIONS
50
51 #POS-TAG dictionary that will be used during the lemmatization process
52 POS_TAG_MAP = defaultdict(lambda : wn.NOUN)
53 POS_TAG_MAP['J'] = wn.ADJ
54 POS_TAG_MAP['V'] = wn.VERB
55 POS_TAG_MAP['R'] = wn.ADV
56
57 #parameters used by the cross validation score function
58 CV_SEED = 7 #seed used for random selection of partitions during cross
59 validation
60 CV_SCORING = 'accuracy'
61
62 #parameters used by the training/set generator
63 TTS_VALIDATION_SIZE = 0.3 #30% of dataset used for training
64 TTS_TEST_SIZE = 1 - TTS_VALIDATION_SIZE #70% of dataset used for testing
65 TTS_SEED = 30 #seed used for random selection of training/test sets
66
67 #parameters used by the gridsearch function
68 GRDSR_SCORING = 'accuracy'
69
70 #DEFINITION OF IMPORTATION, MERGING AND SHUFFLING FUNCTIONS OF DATASETS

```

```

69 def import_dataset(dataset_path, importation_message=None, sep='\t',
70 names=None):
71     """
72     imports a dataset from a given path
73     returns the dataframe containing the imported dataset"""
74
75     print("\n{}".format(importation_message))
76     df = pd.read_csv(dataset_path, sep=sep, header=None, names=names,
77                     encoding='utf-8')
78     print('Size : {}'.format(df.shape))
79     print('Head of imported dataset :')
80     display(df.head())
81
82     return df
83
84 def merge_datasets(df1, df2):
85     """
86     merges datasets contained within dataframes df1 and df2
87     returns a new dataframe containing the merged datasets"""
88
89     df = df1.join(df2)
90
91     print('Size : {}'.format(df.shape))
92     print('Head of merged dataset :')
93     display(df.head())
94
95     return df
96
97 def shuffle_dataset(df):
98     """
99     shuffles dataset entries and reset indexes
100
101     returns a new dataframe containing the shuffled dataset with the
102         reset indexes"""
103
104     shuffled_df = shuffle(df)
105     shuffled_df.reset_index(inplace = True, drop = True)
106
107     print('Head of shuffled dataset :')
108     display(shuffled_df.head())
109
110     return shuffled_df
111
112 #DEFINITION OF PREPROCESSING FUNCTIONS
113 def replace_contractions(document):
114     """
115     replaces contracted expressions in a document
116
117     returns document with no contracted expressions"""
118     return contractions.fix(document)
119
120 def remove_urls(document):

```

```

120     """
121     removes all urls in the document
122
123     returns a document without any urls"""
124     return
125         re.sub(r'https?://(www\.)?[-\w@:~#%]{2,256}\.[a-z]{2,6}\b([-\\w@:~#%_+.~#?&/=
126
127         '', document)
128
129 def remove_empty_html_tags(document):
130     """
131     removes empty html tags like <br />, <hr />, etc.
132
133     returns a document without filtered from empty html tags"""
134     return re.sub(r'(<\w+\s*/?>)', ' ', document)
135
136 def clean_sentence_anchors(document):
137     """
138     cleans all sentences within a document, such that
139     the end of a sentence and the beginning of a new one is separated by
140     a period (or many)
141     followed by a whitespace
142     This cleaning is required because upon removing punctuation,
143     some words get concatenated and create new meaningless terms
144
145     example of a dirty document: "This is a dirty sentence. Another dirty
146     sentence begins"
147     cleaned version: "This is a cleaned sentence. Another cleaned
148     sentence begins"
149
150     This pattern repeats with a sentence ending with a
151     lowercase/uppercase letter and
152     another one beginning with a lowercase/uppercase letter
153     The beginning sentence could also end with a digit and the next
154     sentence could begin with
155     a digit. Hence we get three different patterns:
156     word.*word
157     word.*digit
158     digit.*word
159
160     returns a document with cleaned sentences"""
161
162     word_word = r'([a-zA-Z]+\.[a-zA-Z])\.[a-zA-Z]+' #word(.)word pattern
163     word_digit = r'([a-zA-Z]+\.[a-zA-Z])\.[0-9]+' #word(.)digit pattern
164     digit_word = r'([0-9]+\.[a-zA-Z])\.[a-zA-Z]+' #digit(.)word pattern
165     patterns = [
166         word_word,
167         word_digit,
168         digit_word,
169     ]
170
171     for pattern in patterns:
172         if re.search(pattern, document):
173             document = re.sub(pattern, r'\1. \2', document)

```

```

167     return document
168
169 def remove_non_ascii(tokens):
170     """
171     normalizes the tokens
172     encodes tokens as ASCII characters from tokens
173     and decodes as utf-8
174
175     returns a list of normalized and encoded as ascii tokens"""
176     return [unicodedata.normalize('NFKD', token)
177             .encode('ascii', 'ignore')
178             .decode('utf-8', 'ignore')
179             for token in tokens]
180
181 def split_on_characterset(tokens, regex):
182     """
183     splits a token in tokens upon matching with the character set defined
184     by the regex
185     and appends the tokens obtained from splitting the token to the
186     tokens list
187
188     returns a list of all tokens obtained after splitting problematic
189     tokens"""
190
191     new_tokens = []
192     for token in tokens:
193         if re.search(regex, token) :
194             new_tokens += re.split(regex, token)
195         else:
196             new_tokens.append(token)
197
198     return new_tokens
199
200 def to_lowercase(tokens):
201     """returns a list of tokens in lowercase"""
202     return [token.lower() for token in tokens]
203
204 def replace_numbers(tokens):
205     """
206     replaces tokens representing whole numeric values
207     by their equivalent letter values
208
209     returns a list of transformed tokens"""
210
211     engine = inflect.engine()
212     new_tokens = []
213     for token in tokens:
214         new_token = token
215         if token.isdigit():
216             new_token = engine.number_to_words(token)
217         new_tokens.append(new_token)
218
219     return new_tokens

```

```

218 def remove_punctuation(tokens):
219     """
220     removes tokens not in \w and \s classes of characters.
221     By extension, all punctuation characters will be removed
222
223     returns a list of tokens only in \w and \s"""
224
225     new_tokens = []
226     for token in tokens:
227         new_token = re.sub(r'[\W\s]', '', token)
228         if new_token != '':
229             new_tokens.append(new_token)
230     return new_tokens
231
232 def remove_stopwords(tokens, stopwords=STOP_WORDS):
233     """
234     removes all stopwords (a set) from tokens (a list)
235     except those in exceptions (a set)
236
237     returns a list of tokens that are not stopwords"""
238     return [token for token in tokens if token not in stopwords]
239
240 def lemmatize(tokens, lemmatizer, pos_tag_map):
241     """
242     lematizes all tokens using a lemmatizer and a POS-Tagging map
243
244     returns the list of lemmatized tokens"""
245     return [lemmatizer.lemmatize(token, pos_tag_map[tag[0]]) for token,
246             tag in pos_tag(tokens)]
247
248 def normalize(tokens):
249     """
250     normalizes all tokens by:
251     1. removing non ASCII characters
252     2. converting to lowercase
253     3. splitting wrongfully joined tokens
254     4. replacing numbers with their equivalent letter representation
255     5. removing punctuation
256     6. removing stopwords
257     7. lemmatizing using POS-Tags
258
259     returns the list of normalized tokens"""
260
261     tokens = remove_non_ascii(tokens)
262     tokens = to_lowercase(tokens)
263     tokens = split_on_characterset(tokens, r'[/\~_-]')
264     tokens = replace_numbers(tokens)
265     tokens = remove_punctuation(tokens)
266     tokens = remove_stopwords(tokens)
267     tokens = lemmatize(tokens, WordNetLemmatizer(), POS_TAG_MAP)
268     return tokens
269
270 def preprocess(document):
271     """

```

```

271     preprocesses the document for vectorization by:
272     1. replacing contractions by their equivalent full expressions
273     2. removing empty html tags
274     3. removing urls
275     4. cleaning sentences beginning and end anchors
276     5. tokenizing the document
277     6. normalizing its tokens
278     7. joining normalized tokens back to recreate the document
279
280     returns a preprocessed document, ready for vectorization"""
281
282     document = replace_contractions(document)
283     document = remove_empty_html_tags(document)
284     document = remove_urls(document)
285     document = clean_sentence_anchors(document)
286     tokens = word_tokenize(document)
287     tokens = normalize(tokens)
288     document = ''.join([" " + token for token in tokens]).strip()
289
290     return document
291
292 def preprocess_dataset(dataset):
293     """
294     preprocesses all documents in a dataset
295
296     returns a dataset with preprocessed documents
297     and ready for vectorization"""
298     return [preprocess(document) for document in dataset]
299
300 #DEFINITION OF UTILITY CLASSES
301 #class used to encapsulate the results of the gridsearch
302 class GridSearchResult:
303
304     def __init__(self, name, score, estimator):
305         self.name = name
306         self.score = score
307         self.estimator = estimator
308
309     def __str__(self):
310         return """
311         Model: {}
312         Best Accuracy Score: {}
313         Best Estimator: {}
314         """.format(self.name, self.score, self.estimator)
315
316 #class used to transform a sparse matrix into a dense matrix
317 #to be used by some pipelines during fitting stages
318 class DenseTransformer(TransformerMixin):
319
320     def fit(self, X, y=None, **fit_params):
321         return self
322
323     def transform(self, X, y=None, **fit_params):
324         return X.todense()

```


3 Importation des modules et configuration de l’environnement de travail

Outre les modules essentiels pour faire du machine learning, on importe le fichier `utility_ML.py` contenant l’ensemble des constantes, ressources, fonctions de prétraitement, et classes utilitaires définies.

```
1 #MODULE IMPORTATION AND ENVIRONMENT CONFIGURATION
2 import re
3 import unicodedata
4 import itertools
5 import pickle
6 import warnings
7
8 from time import time
9 from datetime import datetime
10
11 import pandas as pd
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import matplotlib.gridspec as gridspec
15 import contractions
16 import inflect
17
18 from wordcloud import WordCloud, STOPWORDS
19 from collections import defaultdict
20 from mlxtend.plotting import plot_decision_regions
21
22 from nltk import pos_tag
23 from nltk import punkt
24 from nltk.corpus import stopwords, wordnet as wn
25 from nltk.tokenize import word_tokenize
26 from nltk.stem import WordNetLemmatizer
27
28 from sklearn.utils import shuffle
29 from sklearn.base import TransformerMixin
30 from sklearn.feature_extraction.text import TfidfVectorizer
31 from sklearn.model_selection import train_test_split, KFold,
    cross_val_score, GridSearchCV
32 from sklearn.linear_model import LogisticRegression, SGDClassifier
33 from sklearn.neighbors import KNeighborsClassifier
34 from sklearn.tree import DecisionTreeClassifier
35 from sklearn.ensemble import RandomForestClassifier
36 from sklearn.naive_bayes import GaussianNB
37 from sklearn.svm import LinearSVC
38 from sklearn.metrics import accuracy_score, confusion_matrix,
    classification_report
39 from sklearn.pipeline import Pipeline
40
41 from utility_ML import *
```

```

42
43 ##### UNCOMMENT THIS SECTION ON FIRST EXECUTION
44 # import nltk
45 # nltk.download('wordnet')
46 # nltk.download('stopwords')
47 #####
48
49 np.random.seed(500) #set seed for random results base calculation
50 plt.style.use('fivethirtyeight') #choose fivethirtyeight style for plt
51 warnings.filterwarnings('ignore', category=FutureWarning) #filter
    FutureWarnings

```

4 Importation des datasets

On utilise la fonction wrapper `import_dataset()` définie dans le fichier `utility_ML` pour importer les datasets d'avis et de labels `dataset.csv` et `labels.csv`.

Une fois les datasets importés, on les fusionne dans un même dataframe en utilisant la fonction `merge_datasets()` définie dans le fichier `utility_ML`.

Enfin, on mélange les lignes du dataset obtenu d'une manière aléatoire afin de mélanger les avis positifs et négatifs en utilisant la fonction `shuffle_dataset()` définie dans le fichier `utility_ML.py`.

```

1 #IMPORTATION OF DATASETS
2 #importation of opinion dataset
3 df_avis = import_dataset(DATA_PATH, importation_message="\nDataframe des
    avis", sep='\t', names=['Avis'])
4
5 #importation of scores dataset
6 df_score = import_dataset(TARGET_PATH, importation_message='\nDataframe
    des scores', sep='\t', names=['Score'])
7
8 #merging of both datasets
9 df = merge_datasets(df_avis, df_score)
10
11 #shuffling of the merged dataset
12 df = shuffle_dataset(df)

```

Dataframe des avis

Size : (10000, 1)

Head of imported dataset :

```

    Avis
0      Obviously made to show famous 1950s stripper M...
1      This film was more effective in persuading me ...
2      Unless you are already familiar with the pop s...
3      From around the time Europe began fighting Wor...
4      Im not surprised that even cowgirls get the bl...

```

```
Dataframe des scores
Size : (10000, 1)
Head of imported dataset :
```

	Score
0	-1
1	-1
2	-1
3	-1
4	-1

```
Size : (10000, 2)
Head of merged dataset :
```

	Avis	Score
0	Obviously made to show famous 1950s stripper M...	-1
1	This film was more effective in persuading me ...	-1
2	Unless you are already familiar with the pop s...	-1
3	From around the time Europe began fighting Wor...	-1
4	Im not surprised that even cowgirls get the bl...	-1

```
Head of shuffled dataset :
```

	Avis	Score
0	After having read two or three negative review...	1
1	I recently (May 2008) discovered that this chi...	1
2	Pathetic is the word. Bad acting, pathetic scr...	-1
3	Spencer Tracy and Katherine Hepburn would roll...	-1
4	This in my opinion is one of the best action m...	1

5 Pré-traitement des données

On appelle la fonction `preprocess_dataset(dataset)` définie dans le fichier `utility_ML.py` qui appelle la fonction `preprocess(document)` sur la colonne `Avis` afin de prétraiter les documents et les préparer à la vectorisation.

5.1 Pré-traitement d'un document

La fonction `preprocess(document)` effectue les pré-traitements suivants sur chaque document :

1. remplacement des contractions par leurs expressions complètes équivalentes.
2. suppression des balises HTML vides (*e.g.* `
`, `<hr />`, *etc.*)
3. suppression des URLS

4. nettoyage des débuts et fins des phrases (`help clean_sentence_anchors(document)` pour plus d'informations)
5. tokenisation du document
6. normalisation des tokens
7. jointure des tokens normalisés afin de recréer le document, prêt maintenant à la vectorisation

5.2 Pré-traitement des tokens

La fonction `normalize(tokens)` effectue la normalisation des tokens de la manière suivante :

1. normalisation NFKD d'un token et encodage ASCII
2. conversion en minuscule
3. séparation des tokens collés par des caractères de ponctuation (*e.g. token1/token2, token1_token2, etc.*)
4. remplacement des tokens désignant des chiffres par leurs équivalents en lettres
5. suppression des caractères de ponctuation
6. suppression des stopwords (liste des stopwords adaptée au contexte du projet) (*cf. STOP_WORDS dans utility_ML pour plus d'informations*)
7. lemmatization des tokens en utilisant les tags du POS-Tagger

Les fonctions `preprocess_dataset(dataset)`, `preprocess(document)`, `normalize(tokens)` et toutes les fonctions de pré-traitements y inclus sont toutes définies dans `utility_ML.py`.

```
1 #PREPROCESSING DATASET
2 df_transformed = df.copy() #creating a new copy of the dataset that will
   be preprocessed
3 df_transformed['Avis'] = preprocess_dataset(df_transformed['Avis'])
   #preprocessing of opinions column
4 display(df_transformed['Avis'].head())
```

```
0    read two three negative review main page imdb ...
1    recently may two thousand and eight discover c...
2    pathetic word bad act pathetic cheezy dialog h...
3    spencer tracy katherine hepburn would roll gra...
4    opinion one best action 1970s not feature grea...
Name: Avis, dtype: object
```

6 Visualisation des données

Afin de visualiser les effets de pré-traitement on utilisera un **WordCloud** permettant de visualiser les mots les plus fréquents dans les avis positifs et négatifs. Cette visualisation permettra de mieux configurer les fonctions de prétraitement, notamment la liste des stopwords qui pourra être améliorée en y ajoutant les mots neutres propres au domaine désigné par le dataset.

2. Stochastic Gradient Descent Classifier
3. Decision Tree Classifier
4. Random Forest Classifier
5. Gaussian Naive Bayes Classifier
6. K Nearest Neighbors Classifier
7. Linear SVM Classifier (SVM with a linear kernel)

Une fois les modèles testés, on choisira les meilleurs modèles (en termes de score moyen et déviation standard) à calibrer avec un GridSearchCV.

```

1 #CROSS VALIDATION USING ACCURACY METRIC
2 #choosing the data (opinions) and target (score) columns in the dataset
3 X = vectors.toarray()
4 y = df['Score']
5
6 #dictionary containing the models to cross validate using their default
  parameters
7 models = {
8     'LogisticRegression': LogisticRegression(),
9     'SGDClassifier': SGDClassifier(),
10    'DecisionTreeClassifier': DecisionTreeClassifier(),
11    'RandomForestClassifier': RandomForestClassifier(),
12    'GaussianNB': GaussianNB(),
13    'KNeighborsClassifier': KNeighborsClassifier(),
14    'LinearSVC': LinearSVC()
15 }
16 #configuring the parameters used by the cross validation function
17 k_fold = KFold(n_splits=10, shuffle=True, random_state=CV_SEED)
18
19 #cross validation using accuracy metric
20 #for each defined model
21 for name, model in models.items():
22     start_time = time()
23     print('Cross validation started at {}'.format(datetime.now()))
24     cv_score = cross_val_score(model, X, y, cv=k_fold, scoring=CV_SCORING)
25     output = """
26     Time taken to complete cross validation of {}: {} seconds
27     Accuracy scores over 10 evaluations: {}
28     Mean score: {}
29     Standard deviation of scores: {}
30     """.format(name, time() - start_time, cv_score, cv_score.mean(),
31               cv_score.std())
32     print(output)

```

Cross validation started at 2019-04-18 12:38:23.624361

Time taken to complete cross validation of LogisticRegression: 6.682977199554443 seconds
Accuracy scores over 10 evaluations:

[0.914 0.911 0.916 0.922 0.932 0.905 0.933 0.924 0.911 0.921]
Mean score: 0.9189
Standard deviation of scores: 0.008722958213817153

Cross validation started at 2019-04-18 12:39:27.854962

Time taken to complete cross validation of SGDClassifier: 14.247564554214478 seconds
Accuracy scores over 10 evaluations:
[0.9 0.909 0.92 0.934 0.94 0.907 0.936 0.925 0.913 0.925]
Mean score: 0.9209000000000002
Standard deviation of scores: 0.01277849756426787

Cross validation started at 2019-04-18 12:40:48.950553

Time taken to complete cross validation of DecisionTreeClassifier: 796.7198390960693 second
Accuracy scores over 10 evaluations:
[0.757 0.75 0.763 0.755 0.76 0.739 0.755 0.732 0.751 0.753]
Mean score: 0.7515000000000002
Standard deviation of scores: 0.008947066558375441

Cross validation started at 2019-04-18 12:59:27.524247

Time taken to complete cross validation of RandomForestClassifier: 53.930598735809326 second
Accuracy scores over 10 evaluations:
[0.803 0.815 0.813 0.823 0.819 0.833 0.835 0.814 0.799 0.826]
Mean score: 0.818
Standard deviation of scores: 0.011135528725660019

Cross validation started at 2019-04-18 12:56:18.964747

Time taken to complete cross validation of GaussianNB: 23.15500497817993 seconds
Accuracy scores over 10 evaluations:
[0.817 0.841 0.829 0.853 0.836 0.838 0.85 0.847 0.833 0.85]
Mean score: 0.8394
Standard deviation of scores: 0.010650821564555487

Cross validation started at 2019-04-18 19:47:55.137561

Time taken to complete cross validation of KNeighborsClassifier: 2430.809322834015 seconds
Accuracy scores over 10 evaluations:
[0.795 0.793 0.805 0.796 0.799 0.778 0.815 0.775 0.817 0.78]
Mean score: 0.7953
Standard deviation of scores: 0.013849548729110253

Cross validation started at 2019-04-18 12:56:42.120080

Time taken to complete cross validation of LinearSVC: 6.601022243499756 seconds
Accuracy scores over 10 evaluations:


```
[0.913 0.915 0.918 0.926 0.944 0.908 0.938 0.92 0.923 0.925]
Mean score: 0.923
Standard deviation of scores: 0.010497618777608539
```

8.1 Résultats de la cross-validation

En appliquant la cross-validation sur l'ensemble des modèles choisis, on se retrouve avec les résultats suivants ordonnés par ordre décroissant sur les scores moyens :

Modèle	Score moyen	Déviati�n standard
LinearSVC	92%	1%
SGDClassifier	92%	1%
LogisticRegression	91%	0.8%
GaussianNB	84%	1%
RandomForestClassifier	81%	1%
KNeighborsClassifier	79%	1%
DecisionTreeClassifier	75%	0.8%

À partir de ce tableau, nous avons choisi le modèle **Logistic Regression**, ayant un score d'accuracy de 91% et 0.8% de déviation standard et le modèle **Linear SVM** ayant un score d'accuracy de 92% et 1% de déviation standard. On a pu également choisir le modèle **Stochastic Gradient Descent** au lieu du modèle **SVM** ayant un kernel linéaire.

Afin d'exp rimer un peu, nous avons d cider de choisir **Gaussian Naive Bayes** aussi pour visualiser l'effet de choisir des mod les probabilistes dans l'analyse des sentiments.

9 Calibrage des hyperparam tres des mod les en utilisant un Grid-Search

Une fois les mod les choisis   partir de l' tape d' valuation par cross-validation, il faut trouver les meilleurs hyperparam tres permettant de raffiner les r gions de d cision de chaque mod le afin d'avoir les meilleures pr dictions possibles. Ceci est effectu  par un `GridSearchCV`, qui permet de tester diff rentes combinaisons des valeurs des hyperparam tres fournis dans un dictionnaire, pour chaque mod le.

De plus, `GridSearchCV` permet de r p ter le processus sur K partitions, pour choisir les meilleurs r sultats par cross-validation. Nous choisirons ainsi de r p ter le processus sur 5 partitions diff rentes du dataset et nous utiliserons la m trique **Accuracy** pour  valuer les diff rents calibrages des mod les.

Pour le mod le **Logistic Regression**, les hyperparam tres   calibrer sont:

C: la valeur de l'inverse de la r gularization pour la r gression (sauts de 10^{-k} avec $k \in [-4;4]$)

P: la norme   choisir pour les p nalit s (L_1 et L_2)

Pour le mod le **Gaussian Naive Bayes**, il n'existe pas de hyperparam tres   calibrer sauf priors d signant les probabilit s au pr alable estim es pour chacune des classes. Toutefois, nous nous touchons pas   cet hyperparam tre afin qu'il s'adapte dynamiquement au donn es. (cf. https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html)

```

1 #GRIDSEARCH USING THE ACCURACY METRIC FOR PARAMETERS TUNING
2
3 #based on the cross-validation results, using KFold over 10 partitions
4 #the models LogisticRegression and LinearSVC are best suited for the job
5 #However, using GaussianNB should also be taken into account, since it's
6   among the most adapted
7 #for sentiment analysis
8
9 #dictionary containing the candidate models that will be used
10 #for parameters tuning using a GridSearchCV
11 candidates = {
12     'LogisticRegression': models['LogisticRegression'],
13     'LinearSVC': models['LinearSVC']
14 }
15
16 #dictionary of the hyperparameters to be tuned for each model
17 grid_params = {
18     'LogisticRegression': {
19         'C': np.logspace(-4,4,20),
20         'penalty': ['l1','l2']
21     },
22     'LinearSVC': {'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100]}
23 }
24
25 #generation of training/test sets
26 X_train, X_test, y_train, y_test = train_test_split(
27     X,
28     y,
29     train_size = TTS_VALIDATION_SIZE,
30     test_size = TTS_TEST_SIZE,
31     random_state = TTS_SEED
32 )
33
34 #GridSearchCV for every candidate classifier
35 grid_search_results = []
36 for name, model in candidates.items():
37     #creation of the gridsearch
38     grd_sr = GridSearchCV(
39         estimator = model,
40         param_grid = grid_params[name],
41         scoring = GRDSR_SCORING,
42         cv = 5,
43         n_jobs = -1,
44         iid = True,
45         return_train_score = True
46     )
47
48     #execution of the gridsearch
49     start_time = time()
50     print('Grid search started at {}'.format(datetime.now()))
51     grd_sr.fit(X_train, y_train)
52     print('\nTime taken to complete Grid search of {}: {}
53         seconds'.format(name, time() - start_time))

```

```

52     grd_sr_result = GridSearchResult(name, grd_sr.best_score_,
53                                     grd_sr.best_estimator_)
54     print(grd_sr_result)
55     grid_search_results.append(grd_sr_result)
56     #sorting the results by descending order on the score column of the
57     #GridSearchResult objects
58     grid_search_results = sorted(grid_search_results, key=lambda result:
59                                 result.score, reverse=True)
60     print('The best model with the best
61           hyperparameters:\n{}'.format(grid_search_results[0]))

```

Grid search started at 2019-04-19 09:48:26.619439

Time taken to complete Grid search of LogisticRegression: 68.54148316383362 seconds

Model: LogisticRegression

Best Accuracy Score: 0.9073333333333333

Best Estimator: LogisticRegression(C=11.288378916846883, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, max_iter=100, multi_class='warn', n_jobs=None, penalty='l2', random_state=None, solver='warn', tol=0.0001, verbose=0, warm_start=False)

Grid search started at 2019-04-18 20:59:28.571292

Time taken to complete Grid search of LinearSVC: 16.8770534992218 seconds

Model: LinearSVC

Best Accuracy Score: 0.9053333333333333

Best Estimator: LinearSVC(C=1, class_weight=None, dual=True, fit_intercept=True, intercept_scaling=1, loss='squared_hinge', max_iter=1000, multi_class='ovr', penalty='l2', random_state=None, tol=0.0001, verbose=0)

The best model with the best hyperparameters:

Model: LogisticRegression

Best Accuracy Score: 0.9073333333333333

Best Estimator: LogisticRegression(C=11.288378916846883, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, max_iter=100, multi_class='warn', n_jobs=None, penalty='l2', random_state=None, solver='warn', tol=0.0001, verbose=0, warm_start=False)

9.1 Résultats du calibrage avec un GridSearch

En appliquant un GridSearchCV sur l'ensemble des modèles candidats choisis, on se retrouve avec les résultats suivants ordonnés par ordre décroissant sur les scores moyens :

Modèle	Score moyen	Meilleurs Résultats des hyperparamètres
LogisticRegression	90%	$C = 11.288378916846883$; $penalty = L_2$
LinearSVC	90%	$C = 1$

À partir de ce tableau, nous avons choisi le modèle **Logistic Regression**, ayant un score d'accuracy de 90% et les hyperparamètres C et $penalty$ calibrés de la manière suivante : $C = 1$; $penalty = L_2$.

10 Création d'un Pipeline pour le modèle Logistic Regression

Suite au calibrage des hyperparamètres des modèles candidats, on choisit le meilleur modèle calibré avec les meilleurs hyperparamètres pour l'apprentissage et nous créons un pipeline permettant d'enchaîner les étapes de pré-traitement, de vectorisation et d'apprentissage.

Pour la vectorisation, on utilise la technique **BOW (Bag of Words)** avec un `TfidfVectorizer()`. On commence par appliquer les fonctions de prétraitement sur chaque document et effectue la vectorisation pour obtenir la matrice des fréquences des termes, en choisissant les mots ayant une fréquence de document (*i.e. Document Frequency*) de 12, et en appliquant l'algorithme n -grams pour 1 et 2 termes.

Pour l'apprentissage on utilise une instance du modèle choisit avec les meilleurs valeurs trouvées pour les hyperparamètres calibrés, notamment le modèle **Logistic Regression**.

Une fois le pipeline exécuté :

1. on teste le modèle appris
2. on affiche son score de la métrique **Accuracy** en utilisant la fonction `accuracy_score()`
3. on affiche sa matrice de confusion en utilisant la fonction `confusion_matrix()`
4. on affiche les scores des métriques **Precision**, **Recall**, **F1-Score** et **Support** en utilisant la fonction `classification_report()`

Enfin, on sauvegarde le modèle en utilisant la fonction `dump()` du module `pickle`.

```

1 #PIPELINE CREATION FOR LOGISTICREGRESSION CLASSIFIER
2 #creating the pipeline instance
3 lr_pipeline = Pipeline([
4     ('tfidf', TfidfVectorizer(preprocessor=preprocess, min_df=12,
5                             ngram_range=(1, 2))),
6     ('clf', LogisticRegression(C=11.288378916846883, penalty='l2'))
7 ])
8 #choosing data and target columns from initial dataset
9 X = df['Avis']
10 y = df['Score']
11
12 #generating the training/test sets from the initial dataset
13 X_train, X_test, y_train, y_test = train_test_split(
14     X,
15     y,
16     train_size = TTS_VALIDATION_SIZE,
17     test_size = TTS_TEST_SIZE,

```

```

18     random_state = TTS_SEED
19 )
20
21 #learning the model using the pipeline
22 start_time = time()
23 print('LogisticRegression classifier pipeline execution started at
      {}'.format(datetime.now()))
24 lr_pipeline.fit(X_train, y_train)
25 print('\nTime taken to complete pipeline execution: {}
      seconds'.format(time() - start_time))
26
27 #predicting the targets of test data
28 start_time = time()
29 print('\nLogisticRegression classifier prediction started at
      {}'.format(datetime.now()))
30 prediction_result = lr_pipeline.predict(X_test)
31 print('\nTime taken to complete prediction: {} seconds'.format(time() -
      start_time))
32
33 #accuracy, confusion matrix and classification report of the classifier
34 accuracy = accuracy_score(prediction_result, y_test)
35 conf = confusion_matrix(y_test, prediction_result)
36 report = classification_report(y_test, prediction_result)
37 print('')
38 Accuracy: {}
39 Confusion Matrix
40 {}
41
42 Classification Report
43 {}
44 ''.format(accuracy, conf, report))
45
46 #SAVING LOGISTICREGRESSION PIPELINE
47 print('Saving the Logistic Regression pipeline')
48 pickle.dump(lr_pipeline, open(LOGISTICREGRESSION_PATH, 'wb'))

```

LogisticRegression classifier pipeline execution started at 2019-04-19 11:13:34.840659

Time taken to complete pipeline execution: 31.55421018600464 seconds

LogisticRegression classifier prediction started at 2019-04-19 11:14:06.394955

Time taken to complete prediction: 74.41712045669556 seconds

Accuracy: 0.8995714285714286

Confusion Matrix

```

[[3110  373]
 [ 330 3187]]

```

Classification Report

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

-1	0.90	0.89	0.90	3483
1	0.90	0.91	0.90	3517
micro avg	0.90	0.90	0.90	7000
macro avg	0.90	0.90	0.90	7000
weighted avg	0.90	0.90	0.90	7000

Saving the Logistic Regression pipeline

11 Création d'un Pipeline pour le modèle Gaussian Naive Bayes

Nous créons aussi un pipeline pour le modèle **Gaussian Naive Bayes** en utilisant les mêmes paramètres utilisés pour la création du pipeline du modèle **Logistic Regression**. À la différence du modèle **Logistic Regression**, **Gaussian Naive Bayes** travaille avec des **matrices denses** (*i.e. dense matrices*) et non pas des **matrices creuses** (*i.e. sparse matrices*).

Il faut ainsi utiliser un transformateur permettant d'effectuer cette étape suite à l'étape de "fitting" du pipeline. On définit ainsi une classe **DenseTransformer** héritant de la classe de base des transformateurs de matrices du module **scikit-learn** **TransformerMixin** accomplissant cette tâche et fournissant le résultat à la phase d'apprentissage.

Une fois le pipeline exécuté :

1. on teste le modèle appris
2. on affiche son score de la métrique **Accuracy** en utilisant la fonction `accuracy_score()`
3. on affiche sa matrice de confusion en utilisant la fonction `confusion_matrix()`
4. on affiche les scores des métriques **Precision**, **Recall**, **F1-Score** et **Support** en utilisant la fonction `classification_report()`

Enfin, on sauvegarde le modèle en utilisant la fonction `dump()` du module `pickle`.

```

1 #PIPELINE CREATION FOR GAUSSIANNB CLASSIFIER
2 #creating the pipeline instance
3 gnb_pipeline = Pipeline([
4     ('tfidf', TfidfVectorizer(preprocessor=preprocess, min_df=12,
5                             ngram_range=(1, 2))),
6     ('to_dense', DenseTransformer()),
7     ('clf', GaussianNB())
8 ])
9 #choosing data and target columns from initial dataset
10 df_pipeline = df
11 X = df_pipeline['Avis']
12 y = df_pipeline['Score']
13
14 #generating the training/test sets from the initial dataset
15 X_train, X_test, y_train, y_test = train_test_split(
16     X,
17     y,
18     train_size = TTS_VALIDATION_SIZE,
```

```

19     test_size = TTS_TEST_SIZE,
20     random_state = TTS_SEED
21 )
22
23 #learning the model using the pipeline
24 start_time = time()
25 print('GaussianNB classifier pipeline execution started at
      {}'.format(datetime.now()))
26 gnb_pipeline.fit(X_train, y_train)
27 print('\nTime taken to complete pipeline execution: {}
      seconds'.format(time() - start_time))
28
29 #predicting the targets of test data
30 start_time = time()
31 print('\nGaussianNB classifier prediction started at
      {}'.format(datetime.now()))
32 prediction_result = gnb_pipeline.predict(X_test)
33 print('\nTime taken to complete prediction: {} seconds'.format(time() -
      start_time))
34
35 #printing the accuracy, confusion matrix and classification report
36 #of the classifier in the pipeline
37 accuracy = accuracy_score(prediction_result, y_test)
38 conf = confusion_matrix(y_test, prediction_result)
39 report = classification_report(y_test, prediction_result)
40 print('')
41 Accuracy: {}
42 Confusion Matrix
43 {}
44
45 Classification Report
46 {}
47 ''.format(accuracy, conf, report))
48
49 #SAVING GAUSSIANNB PIPELINE
50 print('Saving the Gaussian Naive Bayes pipeline')
51 pickle.dump(gnb_pipeline, open(GAUSSIANNB_PATH, 'wb'))

```

GaussianNB classifier pipeline execution started at 2019-04-19 11:24:59.808241

Time taken to complete pipeline execution: 34.085160970687866 seconds

GaussianNB classifier prediction started at 2019-04-19 11:25:33.893449

Time taken to complete prediction: 73.48156118392944 seconds

Accuracy: 0.8272857142857143

Confusion Matrix

```
[[2814  669]
 [ 540 2977]]
```

Classification Report

	precision	recall	f1-score	support
-1	0.84	0.81	0.82	3483
1	0.82	0.85	0.83	3517
micro avg	0.83	0.83	0.83	7000
macro avg	0.83	0.83	0.83	7000
weighted avg	0.83	0.83	0.83	7000

Saving the Gaussian Naive Bayes pipeline