

1 Classification de données textuelles

Lors de l'étape d'ingénierie de données textuelles nous avons vu que diverses opérations pouvaient être appliquées sur les textes et qu'au final il est possible d'obtenir des textes simplifiés. Nous allons, à présent, étudier comment faire de la classification à partir de données textuelles et comment convertir les textes en vecteurs pour pouvoir faire de la classification.

1.1 Vectorisation

Maintenant qu'un document a été transformé en une séquence de mot, il est nécessaire de la transformer en vecteur. C'est le rôle de la vectorisation.

Bag of Words

La manière la plus simple de vectorisation est d'utiliser les Bag of Words (BOW). Il s'agit, à partir d'une liste de mots (vocabulaire) de compter le nombre d'apparition du mot du vocabulaire dans le document.

Cette opération se fait par :

1. Création d'une instance de la classe CountVectorizer
2. Appel de la fonction fit() pour apprendre le vocabulaire à partir de document
3. Appel de la fonction transform() sur un ou plusieurs documents afin de les encoder dans le vecteur.

Attention, par défaut, CountVectorizer effectue un certain nombre de pré-traitements comme par exemple mise en minuscule. Voir https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html (https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

In []:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2
3 texte = ["This is an example of CountVectorizer for creating a vector"]
4
5 vectorizer = CountVectorizer()
6 # creation du vocabulaire
7 vectorizer.fit(texte)
8 # Contenu du vocabulaire
9 print(vectorizer.vocabulary_)
10 # encodage du document
11 vector = vectorizer.transform(texte)
12
13 print ("Taille du vecteur :\n",vector.shape)
14
```

Il est donc à présent possible de traiter un ensemble de documents comme le montre l'exemple suivant. Nous créons également un dataframe.

In []:

```
1 import pandas as pd
2 from sklearn.feature_extraction.text import CountVectorizer
3 corpus = [
4     'This is my first document.',
5     'This is the document 2 !',
6     'Maybe this is the third document?',
7     'Anything else? may be 40',
8     'Yes !! this is the last one'
9 ]
10 vectorizer = CountVectorizer()
11 X = vectorizer.fit_transform(corpus)
12
13 # vectorizer.get_feature_names()
14 # contient le vocabulaire
15
16
17 df = pd.DataFrame(
18     data=vectorizer.transform(corpus).toarray(),
19     columns=vectorizer.get_feature_names()
20 )
21
22 display(df)
```

Prise en compte des prétraitements

Considérons les prétraitements suivants qui permettent de supprimer les caractères non Ascii, de mettre en minuscule, d'enlever les ponctuations, de remplacer les nombres et d'enlever les stopwords.

In []:

```

1  from nltk.tokenize import word_tokenize
2  import unicodedata
3  import re
4  import inflect
5  from nltk.corpus import stopwords
6
7  def remove_non_ascii(words):
8      new_words = []
9      for word in words:
10         new_word = unicodedata.normalize('NFKD', word).encode('ascii', 'ignore')
11         new_words.append(new_word)
12     return new_words
13
14  def to_lowercase(words):
15      new_words = []
16      for word in words:
17         new_word = word.lower()
18         new_words.append(new_word)
19     return new_words
20
21  def remove_punctuation(words):
22      new_words = []
23      for word in words:
24         new_word = re.sub(r'^\w\s$', '', word)
25         if new_word != '':
26             new_words.append(new_word)
27     return new_words
28
29  def replace_numbers(words):
30      p = inflect.engine()
31      new_words = []
32      for word in words:
33         if word.isdigit():
34             new_word = p.number_to_words(word)
35             new_words.append(new_word)
36         else:
37             new_words.append(word)
38     return new_words
39
40  def remove_stopwords(words):
41      new_words = []
42      for word in words:
43         if word not in stopwords.words('english'):
44             new_words.append(word)
45     return new_words
46
47  def normalize(words):
48      words = remove_non_ascii(words)
49      words = to_lowercase(words)
50      words = replace_numbers(words)
51      words = remove_punctuation(words)
52      words = remove_stopwords(words)
53     return words
54
55  def clean_text(text):
56      tokens = word_tokenize(text)
57      tokens=normalize(tokens)
58      text=" ".join([" "+i for i in tokens]).strip()
59     return text

```

```
60
61
```

In []:

```
1  texte="we have to think that is text is 1000 *#better than. the "
2  print ( "Texte avant le nettoyage \n")
3  print (texte)
4  texte=clean_text(texte)
5  print ( "Texte après le nettoyage \n")
6  print (texte)
```

L'appel aux fonctions de prétraitements peut se faire directement dans CountVectorizer. Attention cependant il est préférable de ne pas le faire. Par exemple dans le cas d'un pipeline et d'un gridsearch le prétraitement sera effectué à chaque fois ! Il est par contre utile de le faire lors de la dernière étape et que le modèle est sauvegardé pour permettre qu'un nouveau document puisse être transformé avant d'être mis sous la forme d'un vecteur (voir plus bas).

In []:

```
1  print (len(corpus))
2  ▼ for i in range(len(corpus)):
3      corpus[i]=clean_text(corpus[i])
```

In []:

```
1
2  #ici le preprocessor ne sert à rien
3  # car les données ont été nettoyées avant.
4  ▼ vectorizer = CountVectorizer(
5      preprocessor=clean_text
6  )
7
8  X = vectorizer.fit_transform(corpus)
9
10
11 ▼ df = pd.DataFrame(
12     data=vectorizer.transform(corpus).toarray(),
13     columns=vectorizer.get_feature_names()
14 )
15
16 display(df)
```

Appel dans CountVectorizer

In []:

```
1
2  ▼ corpus = [
3      'This is my first document.',
4      'This is the document 2 !',
5      'Maybe this is the third document?',
6      'Anything else? may be 40',
7      'Yes !! this is the last one'
8  ]
9  #Rappel ce n'est pas efficace
10 #il vaut mieux traiter les données avant
11 #attention aux pipelines
12  ▼ vectorizer = CountVectorizer(
13      preprocessor=clean_text
14  )
15
16  X = vectorizer.fit_transform(corpus)
17
18
19  ▼ df = pd.DataFrame(
20      data=vectorizer.transform(corpus).toarray(),
21      columns=vectorizer.get_feature_names()
22  )
23
24  display(df)
```

TfidfVectorizer

CountVectorizer en prenant en compte l'occurrence des mots est souvent trop limité. Une alternative est d'utiliser la mesure TF-IDF (Term Frequency – Inverse Document) via une instance de la classe TfidfVectorizer. Le principe est le même que pour CountVectorizer.

Remarque : Il est possible si CountVectorizer a déjà été utilisé de le faire suivre par TfidfTransformer pour simplement mettre à jour les valeurs.

In []:

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 corpus = [
3     'This is my first document.',
4     'This is the document 2 !',
5     'Maybe this is the third document?',
6     'Anything else? may be 40',
7     'Yes !! this is the last one'
8 ]
9 vectorizer = TfidfVectorizer()
10 X = vectorizer.fit_transform(corpus)
11
12 # vectorizer.get_feature_names()
13 # contient le vocabulaire
14
15
16 df = pd.DataFrame(
17     data=vectorizer.transform(corpus).toarray(),
18     columns=vectorizer.get_feature_names()
19 )
20
21 display(df)
```

In []:

```
1 # Appel d'un pré-prétraitement
2 vectorizer = TfidfVectorizer(
3     preprocessor=clean_text
4 )
5
6 X = vectorizer.fit_transform(corpus)
7
8
9 df = pd.DataFrame(
10     data=vectorizer.transform(corpus).toarray(),
11     columns=vectorizer.get_feature_names()
12 )
13
14 display(df)
```

Appel de fonction de prétraitement dans TfidfVectorizer

In []:

```
1  ▾ # Appel d'un pré-prétraitement
2  ▾ vectorizer = TfidfVectorizer(
3      preprocessor=clean_text
4  )
5
6  X = vectorizer.fit_transform(corpus)
7
8
9  ▾ df = pd.DataFrame(
10     data=vectorizer.transform(corpus).toarray(),
11     columns=vectorizer.get_feature_names()
12 )
13
14 display(df)
15
```

Les ngrams

Très souvent il est utile de prendre en compte les n-grammes, i.e. la suite de n mots consécutifs car ils peuvent être important pour la classification. Il est tout à fait possible de les obtenir soit pendant l'étape de prétraitement, soit lors de l'étape de vectorisation en le passant en paramètre.

In []:

```
1  ▾ # Appel d'un pré-prétraitement
2  ▾ vectorizer = TfidfVectorizer(
3      preprocessor=clean_text,
4      ngram_range=(1, 2)
5  )
6
7  X = vectorizer.fit_transform(corpus)
8
9
10 ▾ df = pd.DataFrame(
11     data=vectorizer.transform(corpus).toarray(),
12     columns=vectorizer.get_feature_names()
13 )
14
15 display(df)
```

1.2 Exemple de classification

A partir du moment où nous disposons d'une matrice, nous pouvons appliquer toutes les approches que nous avons vu précédemment.

Nous illustrons au travers d'un exemple de classification multiclasse. Ce dernier est tiré de "The 20 newsgroups text dataset". Il s'agit d'un jeu de données de 20 newsgroups comprenant à peu près 18000 news sur 20 sujets différents.

Ce jeu de données est disponible sous scikit learn qui propose des fonctions pour le manipuler : <https://scikit-learn.org/stable/datasets/index.html#newsgroups-dataset> (<https://scikit-learn.org/stable/datasets/index.html#newsgroups-dataset>)

fetch_20newsgroups permet de charger le fichier. Il est possible de récupérer un jeu d'entraînement, de test ou l'ensemble.

Présentation du jeu de données

In []:

```
1 from sklearn.datasets import fetch_20newsgroups
2 news = fetch_20newsgroups(subset='all')
3 from pprint import pprint
4 print ("liste des topics \n")
5 pprint(list(news.target_names))
```

Téléchargement d'une partie des topics

In []:

```
1 categories = ['alt.atheism', 'talk.religion.misc',
2               'rec.sport.hockey', 'comp.graphics', 'sci.space']
3
4 news = fetch_20newsgroups(subset='all',
5                           categories=categories)
6
7 print ("Taille du jeu de données\n")
8 print (news filenames.shape)
9
10 print ("Un exemple de données\n")
11 print (news.data[5], '\n TOPIC : ', news.target[5],
12        '\n*****')
```

Nettoyage des données

In []:

```
1
2 def clean_news (data):
3     for i in range(len(data)):
4         #print (i)
5         data[i]=clean_text(data[i])
6     return data
7
8
```

In []:

```
1 news.data=clean_news(news.data)
```

Vectorisation

In []:

```
1 vectorizer = TfidfVectorizer()
2 vectors = vectorizer.fit_transform(news.data)
3
4
```

X et y

In []:

```
1  ▼ #Specification des variables à prédire et de la classe X et y
2  X = vectors.toarray()
3
4  y = news.target
5
```

Création du jeu d'apprentissage et de test

In []:

```
1  from sklearn.model_selection import train_test_split
2
3  validation_size=0.3 #30% du jeu de données pour le test
4
5  testsize= 1-validation_size
6  seed=30
7  ▼ X_train,X_test,y_train,y_test=train_test_split(X,
8
9
10                                     y,
11                                     train_size=validation_size,
12                                     random_state=seed,
13                                     test_size=testsize)
```

Utilisation du classifieur

In []:

```
1  from sklearn.naive_bayes import GaussianNB
2  from sklearn.model_selection import KFold
3  from sklearn.model_selection import cross_val_score
4  from time import time
5
6  seed=7
7  k_fold = KFold(n_splits=10, shuffle=True, random_state=seed)
8  clf = GaussianNB()
9
10  scoring = 'accuracy'
11  t0 = time()
12  score = cross_val_score(clf, X, y, cv=k_fold, scoring=scoring)
13  print("Réalisé en %0.3fs" % (time() - t0))
14
15  ▼ print('Les différentes accuracy pour les 10 évaluations sont : \n',
16
17
18          score, '\n')
17  ▼ print ('Accuracy moyenne : ',score.mean(),
18
19          ' standard deviation', score.std())
```

Mise en place d'un pipeline

In []:

```

1  from sklearn.linear_model import SGDClassifier
2  from sklearn.pipeline import Pipeline
3  from sklearn.feature_extraction.text import TfidfTransformer
4  from sklearn.metrics import accuracy_score, confusion_matrix
5  from time import time
6  from sklearn.metrics import classification_report
7
8  ▼ pipeline = Pipeline([('vect', TfidfVectorizer()),
9  ▼                      ('clf', SGDClassifier(loss='hinge',
10                                             penalty='l2',
11                                             alpha=1e-3,
12                                             random_state=42,
13                                             max_iter=5, tol=None)),
14                      ])
15
16
17
18
19  X=news.data
20  y=news.target
21
22
23  ▼ X_train,X_test,y_train,y_test=train_test_split(X,
24                                                  y,
25                                                  train_size=validation_size,
26                                                  random_state=seed,
27                                                  test_size=testsize)
28
29
30  t0 = time()
31  pipeline.fit(X_train, y_train)
32  print("Fit réalisé en %0.3fs" % (time() - t0))
33
34  t0 = time()
35  result = pipeline.predict(X_test)
36  print("Prédiction réalisée en %0.3fs" % (time() - t0))
37
38  print('\n accuracy:',accuracy_score(result, y_test),'\n')
39
40
41
42
43  conf = confusion_matrix(y_test, result)
44  print ('\n matrice de confusion \n',conf)
45
46
47
48  print ('\n',classification_report(y_test, result))
49
50

```

Mise en place d'un gridsearch avec pipeline pour rechercher le meilleur classifieur

Dans cette section nous intégrons un pipeline complet : lancement du TfidfVectorizer , utilisation de deux classifieurs (DecisionTreeClassifier et SGDClassifier) avec les hyperparamètres pour évaluer le meilleur via GridSearchCV. Attention le processus de gridsearch est très long.

In []:

```

1  from sklearn.model_selection import train_test_split
2  from sklearn.model_selection import GridSearchCV
3  from sklearn.pipeline import Pipeline
4  from sklearn.tree import DecisionTreeClassifier
5  from sklearn.linear_model import SGDClassifier
6  from time import time
7  from sklearn.svm import SVC
8  import pickle
9
10
11  # Specification des pipelines
12  # programmation à optimiser par une fonction :)
13  ▼ pipeline_SGDC = Pipeline([('tfidf', TfidfVectorizer()),
14                             ('clf', SGDClassifier())])
15
16
17  ▼ parameters_SGDC = [
18      {'clf__max_iter': (5,),
19       'clf__alpha': (0.00001, 0.000001),
20       'clf__penalty': ('l2', 'elasticnet')}
21  ]
22
23  ▼ pipeline_DT = Pipeline([('tfidf', TfidfVectorizer()),
24                             ('clf', DecisionTreeClassifier())])
25
26
27  #param_range = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
28  param_range = [1, 5, 8, 10]
29  ▼ parameters_DT = [
30      {'clf__min_samples_leaf': param_range,
31       'clf__criterion': ['gini', 'entropy'],
32       'clf__max_depth': param_range,
33       'clf__min_samples_split': param_range[1:]}
34  ]
35
36
37
38
39  X=news.data
40  y=news.target
41
42
43  ▼ X_train,X_test,y_train,y_test=train_test_split(X,
44                                                  y,
45                                                  train_size=validation_size,
46                                                  random_state=seed,
47                                                  test_size=testsize)
48
49  # Creation des GridSearchCV avec les pipelines spécifiques
50
51  ▼ gs_SGDC = GridSearchCV(pipeline_SGDC,
52                           parameters_SGDC,
53                           cv=3,
54                           n_jobs=-1,
55                           scoring='accuracy')
56
57
58  ▼ gs_DT = GridSearchCV(pipeline_DT,

```

```

59         parameters_DT,
60         cv=3,
61         n_jobs=-1,
62         scoring='accuracy')
63
64
65
66 grids = [gs_SGDC, gs_DT]
67 grid_dict={0:'Linear classifiers', 1:'Decision Tree'}
68
69 best_acc = 0.0
70 best_clf = 0.0
71 best_gs = ''
72
73 ▼ for idx,gs in enumerate(grids):
74     print('\nClassifieur: %s' % grid_dict[idx])
75     t0 = time()
76     gs.fit(X_train, y_train)
77     print("Fit réalisé en %0.3fs" % (time() - t0))
78
79     print('Meilleurs paramètres : %s' % gs.best_params_)
80
81     print("Meilleur score d'accuracy sur l'entraînement: %.3f" % gs.best_score_)
82     # Prediction sur le jeu de test avec les meilleurs paramètres
83     t0 = time()
84     result = gs.predict(X_test)
85     print("Prédiction réalisée en %0.3fs" % (time() - t0))
86
87     print("Score d'accuracy pour les meilleurs paramètres sur jeu de test : %")
88
89     print ('\n matrice de confusion \n',confusion_matrix(y_test, result))
90
91     print ('\n',classification_report(y_test, result))
92
93     #Modele avec la meilleure accuracy sur le jeu de test
94 ▼     if accuracy_score(y_test, result) > best_acc:
95         best_acc = accuracy_score(y_test, result)
96         best_gs = gs
97         best_clf = idx
98
99
100
101 ▼ print('\nClassifieur avec la meilleur accuracy sur le jeu de test\n',
102       grid_dict[best_clf])
103
104

```

Recupération du meilleur classifieur avec ses paramètres

Une fois le résultat obtenu, il est possible de récupérer le meilleur classifieur et ses paramètres. Ici il s'agit de Linear Classifier (SGDClassifier). Dans la suite nous montrons comment relancer une classification via un pipeline et une sauvegarde. Attention ici dans le pipeline nous ajoutons le `clean_text` (cf remarques précédentes) et nous sauvegardons le modèle pour pouvoir l'utiliser après avec d'autres données.

Le fait de mettre le `clean_text` dans le pipeline permet lors de la sauvegarde de celui-ci de pouvoir lorsqu'il y a de nouvelles données de les relancer dans le pipeline (les données récupéreront donc la matrice du `TfidfVectorizer` et les pré-traitements associés).

In []:

```

1  from sklearn.linear_model import SGDClassifier
2  from sklearn.pipeline import Pipeline
3  from sklearn.feature_extraction.text import TfidfTransformer
4  from sklearn.metrics import accuracy_score, confusion_matrix
5  from time import time
6  from sklearn.metrics import classification_report
7
8
9  #Recupération des données pour l'exemple
10 #et partir proprement
11 ▼ categories = ['alt.atheism', 'talk.religion.misc',
12                'rec.sport.hockey', 'comp.graphics', 'sci.space']
13
14 ▼ news = fetch_20newsgroups(subset='all',
15                             categories=categories)
16
17
18
19
20
21 ▼ pipeline = Pipeline([('vect', TfidfVectorizer(preprocessor=clean_text)),
22 ▼                  ('clf', SGDClassifier(loss='hinge',
23                                         penalty='l2',
24                                         alpha=1e-05,
25                                         random_state=42,
26                                         max_iter=5,
27                                         tol=None)),
28                        ])
29
30
31
32
33 X=news.data
34 y=news.target
35
36
37 ▼ X_train,X_test,y_train,y_test=train_test_split(X,
38                                                  y,
39                                                  train_size=validation_size,
40                                                  random_state=seed,
41                                                  test_size=testsize)
42
43
44 t0 = time()
45 print ("Lancement du fit \n")
46 pipeline.fit(X_train, y_train)
47 print("Fit réalisé en %0.3fs" % (time() - t0))
48
49 t0 = time()
50 print ("Lancement de la prédiction \n")
51 result = pipeline.predict(X_test)
52 print("Prédiction réalisée en %0.3fs" % (time() - t0))
53
54 print('\n accuracy:',accuracy_score(result, y_test),'\n')
55
56 conf = confusion_matrix(y_test, result)
57 print ('\n matrice de confusion \n',conf)
58
59

```

```

60
61 print ('\n',classification_report(y_test, result))
62
63 print("\nSauvegarde du pipeline grid search")
64 filename = 'thebestone1.pkl'
65 pickle.dump(pipeline, open(filename, 'wb'))
66
67

```

Utilisation de nouvelles données

L'objectif ici est d'utiliser de nouvelles données à partir du modèle appris. Lors de la sauvegarde le pipeline entier a été sauvegardé. Cela implique que lorsque l'on va vouloir prédire les prétraitements du pipeline vont être appliqués.

In []:

```

1 import pickle
2 from sklearn.metrics import accuracy_score
3 from sklearn.metrics import confusion_matrix
4 from sklearn.metrics import classification_report
5
6 print ("Chargement du modèle \n")
7 filename = 'thebestone.pkl'
8 clf_loaded = pickle.load(open(filename, 'rb'))
9
10
11 print ("A partir d'un nouveau texte\n")
12 print ("Utilisation d'un texte de 20newsgroup\n")
13 categories = ['alt.atheism', 'talk.religion.misc',
14               'rec.sport.hockey', 'comp.graphics', 'sci.space']
15
16 news = fetch_20newsgroups(subset='all',
17                           categories=categories)
18
19 print ("Sélection aléatoire de 20 documents \n")
20 from random import randint
21 samples=[]
22 samples_result=[]
23 sample_new=[]
24 for i in range(1,20):
25     val=randint(1,4385)
26     sample_new.append(val)
27     samples.append(news.data[val])
28     samples_result.append(news.target[val])
29
30 print ("Prédiction des news sélectionnées\n")
31
32
33
34 result = clf_loaded.predict(samples)
35
36 print ("Valeurs réelles vs. valeurs prédites\n")
37 for i in range(len(result)):
38     print ("News : ",sample_new[i],
39           "\t réelle ",
40           samples_result[i],
41           " prédite ",
42           result [i])

```

