

Classification avec le jeu de données TITANIC

L'objectif est de poursuivre le TP précédente où de nombreux prétraitements ont été effectués en ingénierie des données pour pouvoir faire un modèle de prédiction des survivants ou non.

Executer les importations suivantes :

In [1]:

```
1  import pandas as pd
2  import numpy as np
3  import sklearn
4  from sklearn.naive_bayes import GaussianNB
5  from sklearn.metrics import accuracy_score
6  from sklearn.model_selection import train_test_split
7  from sklearn.model_selection import KFold
8  from sklearn.model_selection import cross_val_score
9  from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import classification_report
11 import seaborn as sns
12 import matplotlib.pyplot as plt
13 from sklearn.metrics import precision_recall_fscore_support as score
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn.naive_bayes import GaussianNB
18 from sklearn.svm import SVC
19 from sklearn.model_selection import GridSearchCV
20 from sklearn.ensemble import RandomForestClassifier
21
22 #Sickit learn met régulièrement à jour des versions et indique des futurs v
23 #ces deux lignes permettent de ne pas les afficher.
24 import warnings
25 warnings.filterwarnings("ignore", category=FutureWarning)
```

Si vous n'avez pas terminé le TP précédent vous pouvez récupérer le fichier titanic2.csv qui contient les différentes transformations. Le mettre dans le répertoire Dataset.

In [2]:

```
1 df=pd.read_csv('Dataset/titanic2.csv', sep=';')
2 df.head()
```

Out[2]:

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	0	3	1	4	1	0	0	2
1	1	1	0	0	1	0	3	0
2	1	3	0	6	0	0	0	2
3	1	1	0	6	1	0	3	2
4	0	3	1	6	0	0	1	2

Il faut maintenant définir les variables d'apprentissage et la variable à prédire. Ici Survived est la variable à prédire. Attention elle est positionné à la première colonne.

In [3]:

```
1 array = df.values
2 X = array[:,1:8]
3 y = array[:,0]
4
```

Découper le jeu de données en jeu de test et jeu d'apprentissage. Prenez 30% du jeu de données pour le test.

In [4]:

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

Essai d'un classifieur

Avec le classifieur GaussianNB effectuer une première prédiction en donnant la valeur de l'accuracy.

In [5]:

```
1 clf = GaussianNB()
2 clf.fit(X_train, y_train)
3
4 result = clf.predict(X_test)
5
6 ▼ print('\n accuracy :',
7       accuracy_score(result, y_test), '\n')
```

accuracy : 0.7155963302752294

Afficher la matrice de confusion et le classification report.

In [6]:

```
1 conf = confusion_matrix(y_test, result)
2 print ('\n matrice de confusion \n',conf)
3
4 ▼ print ('\n',classification_report(y_test,
5                                     result))
```

matrice de confusion

```
[[57 17]
 [14 21]]
```

	precision	recall	f1-score	support
0	0.80	0.77	0.79	74
1	0.55	0.60	0.58	35
micro avg	0.72	0.72	0.72	109
macro avg	0.68	0.69	0.68	109
weighted avg	0.72	0.72	0.72	109

Effectuer une cross validation avec 10 splits (Kfold)

In [7]:

```
1 seed=7
2 ▼ k_fold = KFold(n_splits=10,
3                 shuffle=True, random_state=seed)
```

Appliquer le classifieur GaussianNB et donner les différentes accuracy pour les 10 évaluations.

In [8]:

```
1 clf = GaussianNB()
2
3 scoring = 'accuracy'
4 ▼ score = cross_val_score(clf, X, y,
5                             cv=k_fold, scoring=scoring)
6
7 print('Les différentes accuracy pour les 10 évaluations sont : \n',score,'')
8 ▼ print ('Accuracy moyenne : ',score.mean(),
9         ' standard deviation', score.std())
10
```

Les différentes accuracy pour les 10 évaluations sont :

```
[0.8125      0.8125      0.75      0.75      0.875      0.86666667
 0.73333333 0.8        0.73333333 0.73333333]
```

Accuracy moyenne : 0.7866666666666667 standard deviation 0.051901348730066764

Essai de plusieurs classifieurs

Utiliser à présent différents classifieurs : KNeighborsClassifier, DecisionTreeClassifier, GaussianNB, SVC et RandomForestClassifier.

In [9]:

```
1 seed = 7
2 scoring = 'accuracy'
3 models = []
4 models.append(('KNN', KNeighborsClassifier()))
5 models.append(('CART', DecisionTreeClassifier()))
6 models.append(('NB', GaussianNB()))
7 models.append(('SVM', SVC(gamma='auto')))
8 models.append(('RFO', RandomForestClassifier()))
```

Quels sont les résultats des classifications ? Quel est le classifieur le plus efficace ?

In [10]:

```
1 results = []
2 names = []
3 ▼ for name,model in models:
4     kfold = KFold(n_splits=10, random_state=seed)
5 ▼     cv_results = cross_val_score(model, X,
6                                     y, cv=kfold, scoring=scoring)
7     results.append(cv_results)
8     names.append(name)
9 ▼     msg = "%s: %f (%f)" % (name, cv_results.mean(),
10                             cv_results.std())
11     print(msg)
```

```
KNN: 0.647083 (0.109117)
CART: 0.728333 (0.050649)
NB: 0.768750 (0.072821)
SVM: 0.730000 (0.071317)
RFO: 0.774167 (0.099055)
```

A l'aide de boxplot afficher les résultats des différents classifieurs.

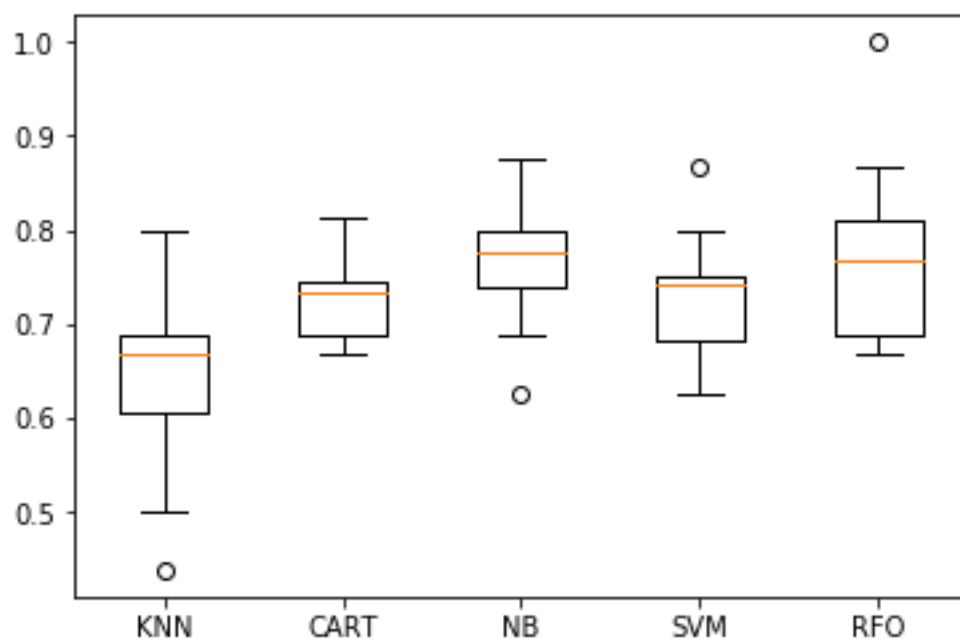
In [11]:

```
1 fig = plt.figure()
2 fig.suptitle('Comparaison des algorithmes')
3 ax = fig.add_subplot(111)
4 plt.boxplot(results)
5 ax.set_xticklabels(names)
```

Out[11]:

```
[Text(0, 0, 'KNN'),
Text(0, 0, 'CART'),
Text(0, 0, 'NB'),
Text(0, 0, 'SVM'),
Text(0, 0, 'RFO')]
```

Comparaison des algorithmes



Le classifieur RandomForestClassifier donne de meilleurs résultats. A l'aide de GridSearch, évaluer les différents paramètres :

```
grid_param = {'n_estimators': [4, 6, 9],  
'max_features': ['log2', 'sqrt', 'auto'],  
'criterion': ['entropy', 'gini'],  
'max_depth': [2, 3, 5, 10],  
'min_samples_split': [2, 3, 5],  
'min_samples_leaf': [1, 5, 8]  
}
```

In [12]:

```
1  ▼ grid_param = {'n_estimators': [4, 6, 9],  
2                  'max_features': ['log2', 'sqrt', 'auto'],  
3                  'criterion': ['entropy', 'gini'],  
4                  'max_depth': [2, 3, 5, 10],  
5                  'min_samples_split': [2, 3, 5],  
6                  'min_samples_leaf': [1, 5, 8]  
7                  }  
8  
9  ▼ gd_sr = GridSearchCV(estimator=RandomForestClassifier(),  
10                        param_grid=grid_param,  
11                        scoring='accuracy',  
12                        cv=5,  
13                        n_jobs=-1,  
14                        iid=True,  
15                        return_train_score=True)  
16  
17
```

Quelle est la meilleure accuracy et pour quels paramètres ?

In [13]:

```
1 gd_sr.fit(X_train, y_train)
2 ▼ print ('meilleur score ',
3         gd_sr.best_score_, '\n')
4 ▼ print ('meilleurs paramètres',
5         gd_sr.best_params_, '\n')
6 ▼ print ('meilleur estimateur',
7         gd_sr.best_estimator_, '\n')
```

meilleur score 0.782608695652174

meilleurs paramètres {'criterion': 'entropy', 'max_depth': 2, 'max_features': 'sqrt', 'min_samples_leaf': 5, 'min_samples_split': 2, 'n_estimators': 4}

meilleur estimateur RandomForestClassifier(bootstrap=True, class_weight=None, criterion='entropy',
max_depth=2, max_features='sqrt', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=5, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=4, n_jobs=None,
oob_score=False, random_state=None, verbose=0,
warm_start=False)

GaussianNB obtient un assez bon score mais il n'a pas d'hyperparamètres à rechercher. Il est suivi par DecisionTreeClassifier et SVC.

Appliquer GridSearchCV pour ces deux classifieurs.

In [14]:

```
1  from sklearn.tree import DecisionTreeClassifier
2
3
4  ▼ grid_param = {
5      'max_depth': [1,2,3,4,5,6,7,8,9,10],
6      'criterion': ['gini', 'entropy'],
7      'min_samples_leaf': [1,2,3,4,5,6,7,8,9,10]
8  }
9
10
11  ▼ gd_sr = GridSearchCV(estimator=DecisionTreeClassifier(),
12                        param_grid=grid_param,
13                        scoring='accuracy',
14                        cv=10,
15                        n_jobs=-1,
16                        iid=True,
17                        return_train_score=True)
18
19  gd_sr.fit(X_train, y_train)
20
21  gd_sr.fit(X_train, y_train)
22  print ('meilleur score ',gd_sr.best_score_,'\n')
23  print ('meilleurs paramètres', gd_sr.best_params_,'\n')
24  print ('meilleur estimateur',gd_sr.best_estimator_,'\n')
```

meilleur score 0.782608695652174

meilleurs paramètres {'criterion': 'gini', 'max_depth': 1, 'min_samples_leaf': 1}

meilleur estimateur DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=1, max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, presort=False, random_state=None, splitter='best')

In [15]:

```
1  ▼ grid_param = {
2      'C': [0.001, 0.01, 0.1, 1, 10],
3      'gamma': [0.001, 0.01, 0.1, 1],
4      'kernel': ['linear', 'rbf']}
5
6  ▼ gd_sr = GridSearchCV(estimator=SVC(),
7                        param_grid=grid_param,
8                        scoring='accuracy',
9                        cv=5,
10                       n_jobs=1,
11                       iid=True,
12                       return_train_score=True)
13
14  gd_sr.fit(X_train, y_train)
15
16  ▼ print ('meilleur score ',
17         gd_sr.best_score_, '\n')
18  ▼ print ('meilleurs paramètres',
19         gd_sr.best_params_, '\n')
20  ▼ print ('meilleur estimateur',
21         gd_sr.best_estimator_, '\n')
22
23
```

meilleur score 0.7391304347826086

meilleurs paramètres {'C': 1, 'gamma': 0.001, 'kernel': 'linear'}

meilleur estimateur SVC(C=1, cache_size=200, class_weight=None, coef
0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.001, kernel='line
ar',
max_iter=-1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False)

Vous constatez que RandomForestClassifier et DecisionTreeClassifier peuvent avoir la même accuracy. Faire un gridsearch avec les deux classifieurs en prenant les paramètres précédents.

In [16]:

```
1  ▼ classifiers = {
2      'RandomForestClassifier': RandomForestClassifier(),
3      'DecisionTreeClassifier': DecisionTreeClassifier()
4  }
5
6  ▼ params = {'RandomForestClassifier' :
7  ▼          [{'n_estimators': [4, 6, 9]],
8              {'max_features': ['log2', 'sqrt', 'auto']},
9              {'criterion': ['entropy', 'gini']}},
10             {'max_depth': [2, 3, 5, 10]},
11             {'min_samples_split': [2, 3, 5]},
12             {'min_samples_leaf': [1, 5, 8]}}],
13             'DecisionTreeClassifier':
14             [{'max_depth': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]},
15             {'criterion': ['gini', 'entropy']},
16             {'min_samples_leaf': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]}}]}
17
18
```

In [17]:

```
1  ▼ class Result:
2  ▼     def __init__(self, name, score, parameters):
3          self.name = name
4          self.score = score
5          self.parameters = parameters
6  ▼     def __repr__(self):
7          return repr((self.name, self.score, self.parameters))
8
9
10 results = []
11 ▼ for key, value in classifiers.items():
12 ▼     gd_sr = GridSearchCV(estimator=value,
13                           param_grid=params[key],
14                           scoring='accuracy',
15                           cv=10,
16                           n_jobs=1,
17                           iid=True)
18     gd_sr.fit(X_train, y_train)
19 ▼     result=Result(key, gd_sr.best_score_,
20                   gd_sr.best_estimator_)
21     results.append(result)
22
23
24
25 ▼ results=sorted(results,
26                  key=lambda result: result.score,
27                  reverse=True)
28
29 print ('Le meilleur resultat : \n')
30 ▼ print ('Classifier : ', results[0].name,
31         ' score %.4f' % results[0].score,
32         ' avec ', results[0].parameters, '\n')
33
34 print ('Tous les résultats : \n')
```

```

35  ▼ for result in results:
36  ▼     print ('Classifier : ',result.name,
37             ' score %0.4f' %result.score,
38             ' avec ',result.parameters,'\n')
39
40
41

```

Le meilleur resultat :

```

Classifier : DecisionTreeClassifier score 0.7826 avec DecisionTr
eeClassifier(class_weight=None, criterion='gini', max_depth=1,
             max_features=None, max_leaf_nodes=None,
             min_impurity_decrease=0.0, min_impurity_split=None,
             min_samples_leaf=1, min_samples_split=2,
             min_weight_fraction_leaf=0.0, presort=False, random_stat
e=None,
             splitter='best')

```

Tous les résultats :

```

Classifier : DecisionTreeClassifier score 0.7826 avec DecisionTr
eeClassifier(class_weight=None, criterion='gini', max_depth=1,
             max_features=None, max_leaf_nodes=None,
             min_impurity_decrease=0.0, min_impurity_split=None,
             min_samples_leaf=1, min_samples_split=2,
             min_weight_fraction_leaf=0.0, presort=False, random_stat
e=None,
             splitter='best')

```

```

Classifier : RandomForestClassifier score 0.6957 avec RandomFore
stClassifier(bootstrap=True, class_weight=None, criterion='gini',
             max_depth=None, max_features='auto', max_leaf_nodes=None
,
             min_impurity_decrease=0.0, min_impurity_split=None,
             min_samples_leaf=1, min_samples_split=5,
             min_weight_fraction_leaf=0.0, n_estimators=10, n_jobs=No
ne,
             oob_score=False, random_state=None, verbose=0,
             warm_start=False)

```

Conclusion : il faut toujours faire attention aux paramètres passés dans un classifieur. Ils peuvent tout changer !

Sauvegarde du modèle

Sauvegarder le meilleur modèle appris et recharger le pour le tester avec `y_test`.

In [18]:

```
1 clf=gd_sr.best_estimator_  
2 import pickle  
3 filename = 'pkl_Titanicbestmodel.sav'  
4 pickle.dump(clf, open(filename, 'wb'))
```

In [19]:

```
1 clf_loaded = pickle.load(open(filename, 'rb'))  
2 print ('Modèle chargé',clf_loaded,'\n')  
3 result = clf_loaded.predict(X_test)  
4  
5 print('\n accuracy:\n')  
6 print (accuracy_score(result, y_test),'\n')  
7  
8 conf = confusion_matrix(y_test, result)  
9 print ('\n matrice de confusion \n',conf)  
10 print ('\n',classification_report(y_test, result))
```

```
Modèle chargé DecisionTreeClassifier(class_weight=None, criterion='g  
ini', max_depth=1,  
    max_features=None, max_leaf_nodes=None,  
    min_impurity_decrease=0.0, min_impurity_split=None,  
    min_samples_leaf=1, min_samples_split=2,  
    min_weight_fraction_leaf=0.0, presort=False, random_stat  
e=None,  
    splitter='best')
```

accuracy:

0.8165137614678899

matrice de confusion

```
[[62 12]  
 [ 8 27]]
```

	precision	recall	f1-score	support
0	0.89	0.84	0.86	74
1	0.69	0.77	0.73	35
micro avg	0.82	0.82	0.82	109
macro avg	0.79	0.80	0.80	109
weighted avg	0.82	0.82	0.82	109

