

Les Pima sont un groupe d'Amérindiens vivant en Arizona. Une prédisposition génétique a permis à ce groupe de survivre normalement à un régime pauvre en glucides pendant des années. Au cours des dernières années, en raison d'un passage soudain des cultures agricoles traditionnelles aux aliments transformés, ainsi que d'un déclin de l'activité physique, ils ont développé la prévalence la plus élevée de diabète de type 2 et pour cette raison ils ont fait l'objet de nombreuses études.

L'ensemble de données comprend des données provenant de 768 femmes présentant 8 caractéristiques, en particulier :

- Nombre de grossesses
- Concentration de glucose plasmatique a 2 heures dans un test de tolérance au glucose par voie orale
- Tension artérielle diastolique (mm Hg)
- Epaisseur du pli cutané du triceps (mm)
- 2 heures d'insuline sérique (mu U/ml)
- Indice de masse corporelle (poids en kg/(taille en m)^2)
- Le diabète pedigree fonction
- Âge (années) La dernière colonne de l'ensemble de données indique si la personne a reçu un diagnostic de diabète (1) ou non (0).

L'objectif est de déterminer quelles sont les caractéristiques (features) pour identifier les personnes qui ont un diabète de type 2

Récupérer le fichier pima-indians-diabetes.csv et le mettre dans un dataframe. Attention les premières lignes correspondent à la description des données. Il est possible de ne pas les lire en mettant skiprows=9 dans la fonction read_csv.

In [56]:

```
1 import pandas as pd
2 import numpy as np
3 names=[
4     "NumTimesPrg", "PlGlcConc", "BloodP",
5     "SkinThick", "TwoHourSerIns", "BMI",
6     "DiPedFunc", "Age", "HasDiabetes"]
7 #il faut sauter les 9 premières lignes qui sont le descriptif des variables
8 df = pd.read_csv('pima-indians-diabetes.csv', names=names, skiprows=9)
9
```

Afficher le nombre de ligne et de colonnes du dataframe ainsi que les 5 premières lignes

In [8]:

```
1 print (df.shape)
2 display(df.head())
```

(768, 9)

	NumTimesPrg	PIGlcConc	BloodP	SkinThick	TwoHourSerIns	BMI	DiPedFunc	Age	HasD
0	6	148	72	35	0	33.6	0.627	50	
1	1	85	66	29	0	26.6	0.351	31	
2	8	183	64	0	0	23.3	0.672	32	
3	1	89	66	23	94	28.1	0.167	21	
4	0	137	40	35	168	43.1	2.288	33	

Afficher la matrice de corrélation. Rappel il faut utiliser la fonction corr().

In [10]:

```
1 corr = df.corr()
2 display(corr)
```

	NumTimesPrg	PIGlcConc	BloodP	SkinThick	TwoHourSerIns	BMI	DiF
NumTimesPrg	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.033523
PIGlcConc	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337
BloodP	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265
SkinThick	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928
TwoHourSerIns	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071
BMI	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647
DiPedFunc	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000
Age	0.544341	0.263514	0.239528	-0.113970	-0.042163	0.036242	0.036242
HasDiabetes	0.221898	0.466581	0.065068	0.074752	0.130548	0.292695	0.292695

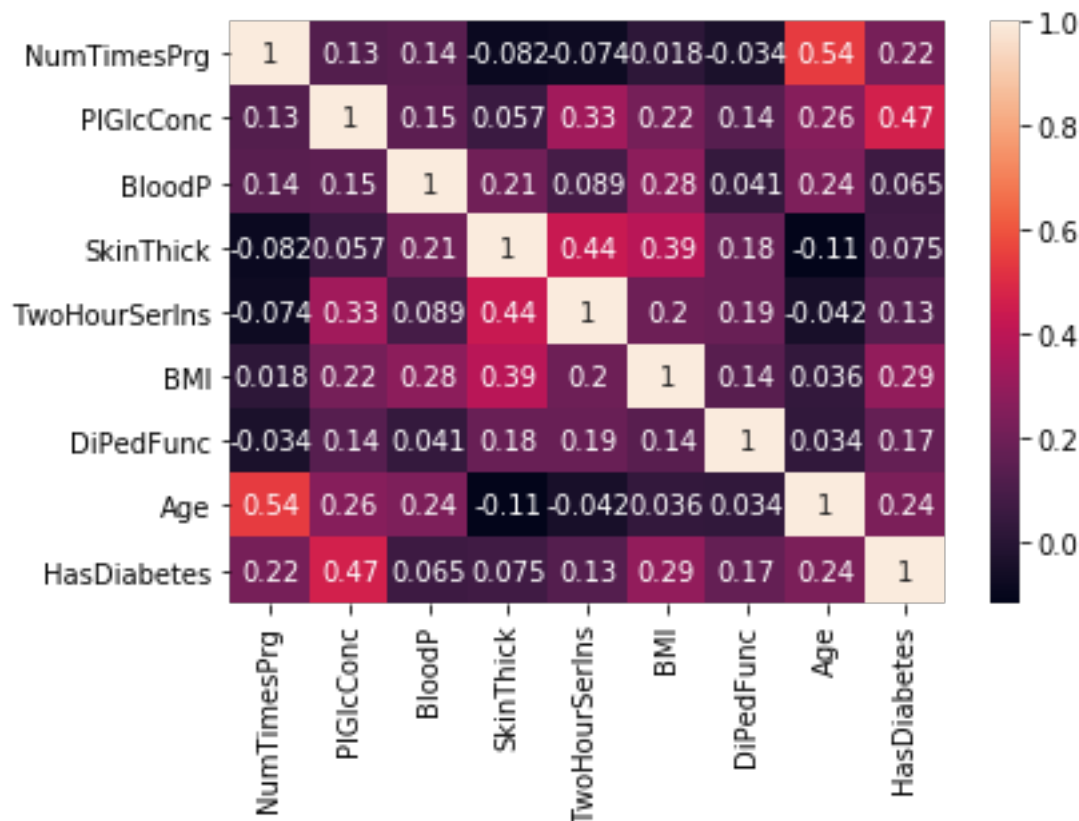
Afficher, à l'aide de seaborn, la matrice de correlation

In [12]:

```
1 %matplotlib inline
2 import seaborn as sns
3 sns.heatmap(corr, annot = True)
```

Out[12]:

<matplotlib.axes._subplots.AxesSubplot at 0x11428a518>



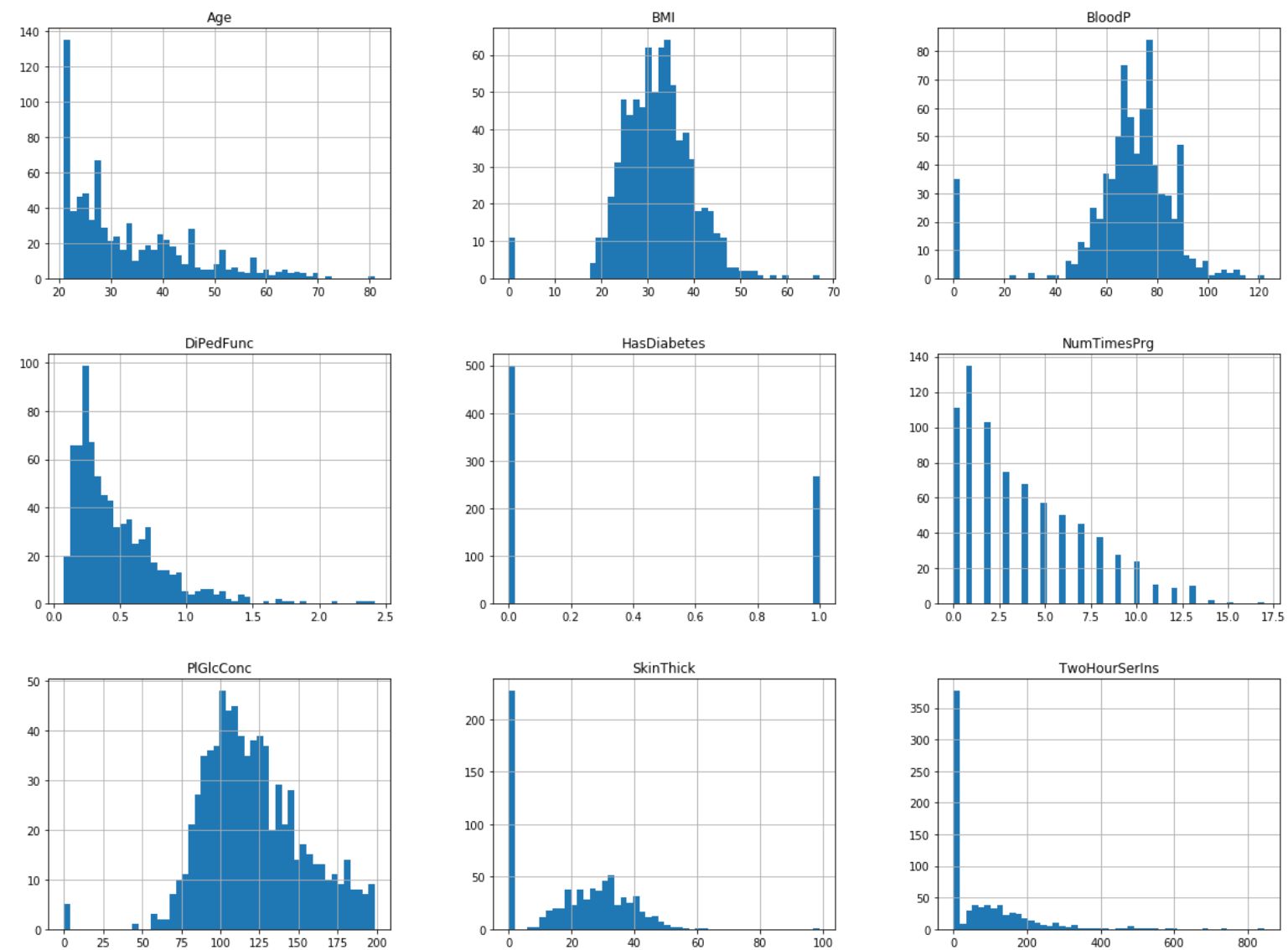
Il est important d'analyser les histogrammes de chaque variable pour mieux comprendre comment les données sont réparties.

A l'aide du code suivant, afficher les différents histogrammes.

```
import matplotlib.pyplot as plt
df.hist(bins=50, figsize=(20, 15))
plt.show()
```

In [15]:

```
1 import matplotlib.pyplot as plt
2 df.hist(bins=50, figsize=(20, 15))
3 plt.show()
```



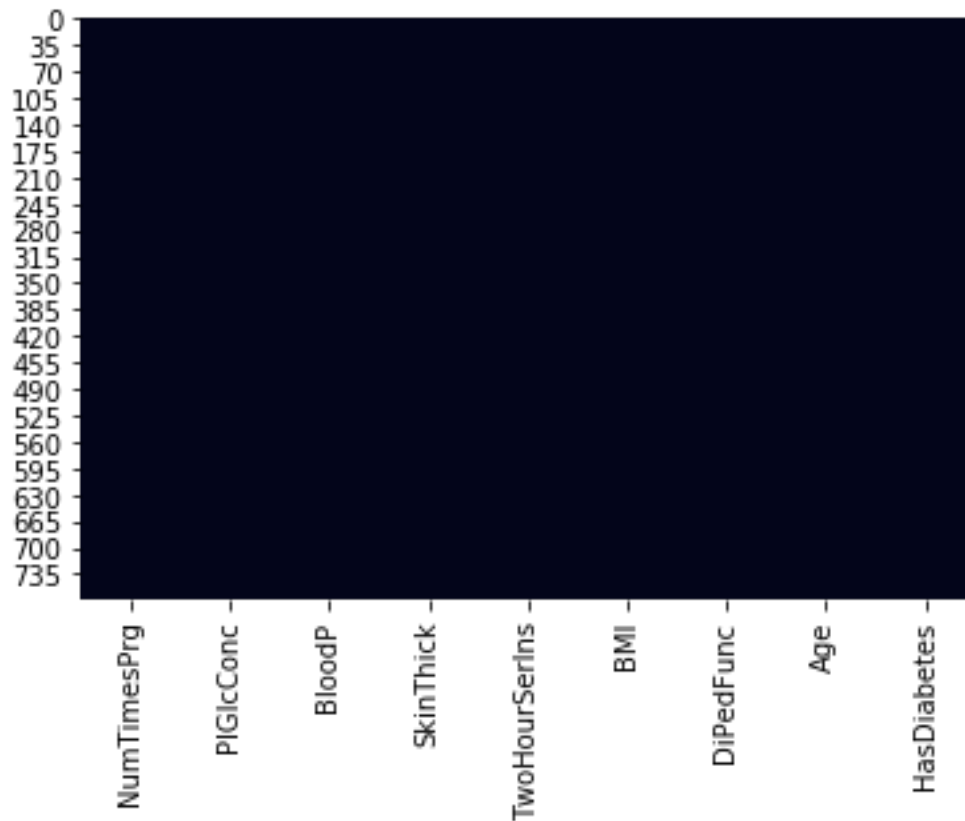
Existe-t'il des valeurs nulles ? Existe-t-il des valeurs manquantes ? Rappel vous pouvez le voir avec des histogrammes mais aussi avec une heatmap.

In [57]:

```
1 import seaborn as sns
2 sns.heatmap(df.isnull(), cbar=False)
```

Out[57]:

<matplotlib.axes._subplots.AxesSubplot at 0x117ad12e8>



En fait on peut constater qu'il n'y a pas de valeurs manquantes avec le heatmap mais par contre il y a des valeurs nulles. Il faut toujours faire attention à la manière dont sont codées les valeurs manquantes. Ici nous voyons dans les histogrammes que pour BMI, BloodP, PIGlcConc, SkinThick, TwoHourSerIns il existe des valeurs manquantes. Le nombre de grossesses n'est pas considéré comme une valeur manquante bien sûr.

Transformer les valeurs nulles par la médiane de la série.

In [18]:

```
1  ▼  # valeur médiane
2      median_bmi = df['BMI'].median()
3      # remplacement par la médiane
4  ▼  df['BMI'] = df['BMI'].replace(
5      to_replace=0, value=median_bmi)
6
7      # valeur médiane
8      median_bloodP = df['BloodP'].median()
9      # remplacement par la médiane
10 ▼  df['BloodP'] = df['BloodP'].replace(
11     to_replace=0, value=median_bloodP)
12
13     # valeur médiane
14     median_PlGlcConc = df['PlGlcConc'].median()
15     # remplacement par la médiane
16 ▼  df['PlGlcConc'] = df['PlGlcConc'].replace(
17     to_replace=0, value=median_PlGlcConc)
18
19     # valeur médiane
20     median_SkinThick = df['SkinThick'].median()
21     # remplacement par la médiane
22 ▼  df['SkinThick'] = df['SkinThick'].replace(
23     to_replace=0, value=median_SkinThick)
24
25     # valeur médiane
26     median_TwoHourSerIns = df['TwoHourSerIns'].median()
27     # remplacement par la médiane
28 ▼  df['TwoHourSerIns'] = df['TwoHourSerIns'].replace(
29     to_replace=0, value=median_TwoHourSerIns)
30
```

Les données sont, à présent, transformées et nous allons pouvoir créer un jeu de données de test et d'apprentissage. Faire une copie du dataframe en df2. Sur df appliquer un scaling pour normaliser les valeurs par rapport à la moyenne et l'écart type (utilisation de StandardScaler()). Nous conservons la copie df2 sans transformation.

L'objectif à présent est d'appliquer différents classifieurs pour voir celui qui est le plus performant. Pour le ou les meilleurs rechercher les hyperparamètres et créer un pipeline à sauvegarder. Il faut ensuite pouvoir traiter de nouvelles données pour prédire si il y a diabète ou pas.

Tester les résultats sur df et sur df2.

[illegible]

	0	1	2	3	4	5
6 \						
0	0.639947	0.848324	0.149641	0.907270	-0.692891	0.204013
492						
1	-0.844885	-1.123396	-0.160546	0.530902	-0.692891	-0.684422
061						
2	1.233880	1.943724	-0.263941	-1.288212	-0.692891	-1.103255
397						
3	-0.844885	-0.998208	-0.160546	0.154533	0.123302	-0.494043
763						
4	-1.141852	0.504055	-1.504687	0.907270	0.765836	1.409746
909						
	7					
0	1.425995					
1	-0.190672					
2	-0.105584					
3	-1.041549					
4	-0.020496					

[illegible]

	0	1	2	3	4	5
6 \						
0 492	0.639947	0.848324	0.149641	0.907270	-0.692891	0.204013
1 061	-0.844885	-1.123396	-0.160546	0.530902	-0.692891	-0.684422
2 397	1.233880	1.943724	-0.263941	-1.288212	-0.692891	-1.103255
3 763	-0.844885	-0.998208	-0.160546	0.154533	0.123302	-0.494043
4 909	-1.141852	0.504055	-1.504687	0.907270	0.765836	1.409746

```

7
0  1.425995
1 -0.190672
2 -0.105584
3 -1.041549
4 -0.020496

```


In [71]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.svm import SVC
5 from sklearn.svm import LinearSVC
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.tree import DecisionTreeRegressor
8 from sklearn import model_selection
```

In [72]:

```
1 models = []
2 models.append(('LR', LogisticRegression()))
3 models.append(('KNN', KNeighborsClassifier()))
4 models.append(('NB', GaussianNB()))
5 models.append(('SVC', SVC(gamma='auto')))
6 models.append(('LSVC', LinearSVC(max_iter=3000)))
7 models.append(('RFC', RandomForestClassifier()))
8 models.append(('DTR', DecisionTreeRegressor()))
```

In [73]:

```
1  import warnings
2  warnings.filterwarnings("ignore", category=FutureWarning)
3
4  from sklearn.metrics import confusion_matrix
5  from sklearn.metrics import classification_report
6  from sklearn.model_selection import KFold
7  from sklearn.model_selection import cross_val_score
8  import time
9  seed = 7
10 results = []
11 names = []
12 scoring='accuracy'
13 ▼ for name,model in models:
14     kfold = KFold(n_splits=10, random_state=seed)
15     start_time = time.time()
16     cv_results = cross_val_score(model, X, y, cv=kfold, scoring=scoring)
17     print ("Time pour",name," ",time.time() - start_time)
18     results.append(cv_results)
19     names.append(name)
20     msg = "%s: %f (%f)" % (name, cv_results.mean(), cv_results.std())
21     print(msg)
22
```

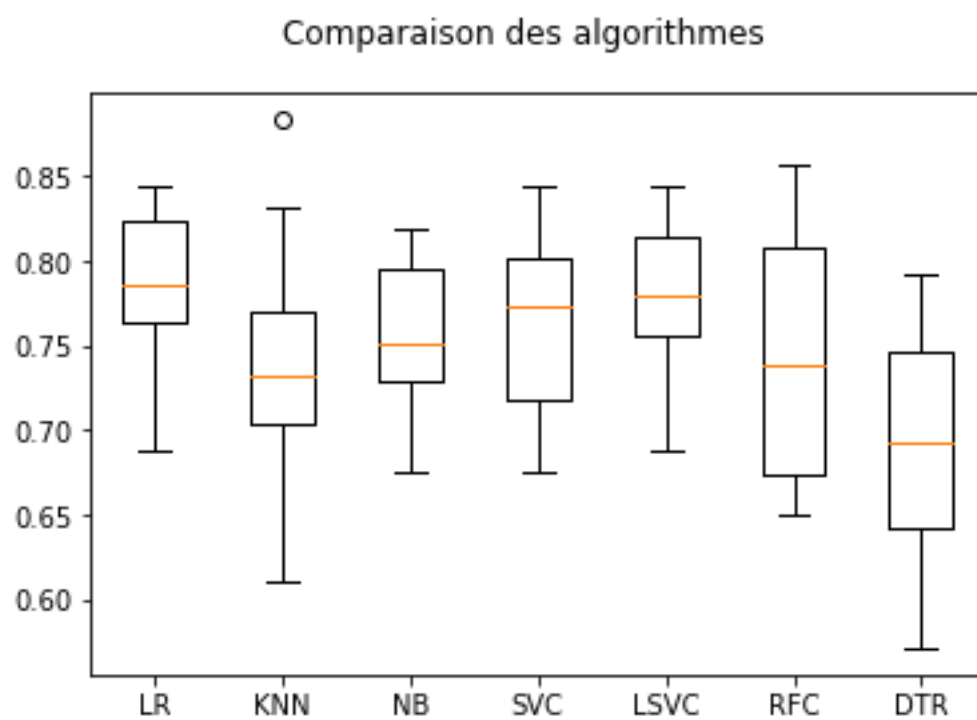
```
Time pour LR    0.048712968826293945
LR: 0.779956 (0.050088)
Time pour KNN   0.030850887298583984
KNN: 0.742139 (0.071500)
Time pour NB    0.015557050704956055
NB: 0.755178 (0.042766)
Time pour SVC   0.12708377838134766
SVC: 0.764286 (0.056962)
Time pour LSVC  0.4412078857421875
LSVC: 0.776059 (0.049143)
Time pour RFC   0.1575641632080078
RFC: 0.743438 (0.072740)
Time pour DTR   0.028022050857543945
DTR: 0.692618 (0.066232)
```

In [74]:

```
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[74]:

```
[Text(0, 0, 'LR'),
 Text(0, 0, 'KNN'),
 Text(0, 0, 'NB'),
 Text(0, 0, 'SVC'),
 Text(0, 0, 'LSVC'),
 Text(0, 0, 'RFC'),
 Text(0, 0, 'DTR')]
```



Les meilleurs algorithmes sont LR et SVC. Nous pouvons rechercher les hyperparamètres pour ces deux algorithmes.

In [75]:

```
1  from sklearn.model_selection import GridSearchCV
2
3  ▼ grid_param = {
4      'C' : [0.001,0.01,0.1,1,10,100]
5  }
6
7
8  ▼ gd_sr = GridSearchCV(estimator=LogisticRegression(),
9                        param_grid=grid_param,
10                       scoring='accuracy',
11                       cv=5,
12                       n_jobs=-1,
13                       iid=True,
14                       return_train_score=True)
15
16  gd_sr.fit(X_train, y_train)
17
18  print ('meilleur score ',gd_sr.best_score_,'\n')
19  print ('meilleurs paramètres', gd_sr.best_params_,'\n')
20  print ('meilleur estimateur',gd_sr.best_estimator_,'\n')
```

meilleur score 0.7516339869281046

meilleurs paramètres {'C': 1}

meilleur estimateur LogisticRegression(C=1, 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)

In [52]:

```
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 ',gd_sr.best_score_,'\n')
17  print ('meilleurs paramètres', gd_sr.best_params_,'\n')
18  print ('meilleur estimateur',gd_sr.best_estimator_,'\n')
```

meilleur score 0.7608695652173914

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)

SVC s'avère obtenir de meilleurs résultats. Nous l'utilisons pour faire de la prédiction sur les données test.

In [76]:

```
1  from sklearn.metrics import accuracy_score
2  from sklearn.metrics import confusion_matrix
3  from sklearn.metrics import classification_report
4  # Creation d'une instance de l'algorithme en utilisant les meilleurs param
5  svc = gd_sr.best_estimator_
6
7  validation_size=0.2 #20% du jeu de données pour le test
8
9  testsize= 1-validation_size
10 seed=20
11 X_train,X_test,y_train,y_test=train_test_split(X, y,
12                                                train_size=validation_size,
13                                                random_state=seed,
14                                                test_size=testsize)
15
16
17 svc.fit(X_train, y_train)
18 result = svc.predict(X_test)
19 print('\n accuracy: ', accuracy_score(result, y_test),'\n')
20
21 conf = confusion_matrix(y_test, result)
22 print ('\n matrice de confusion \n',conf)
23 print ('\n',classification_report(y_test, result))
```

accuracy: 0.7430894308943089

matrice de confusion

```
[[363  26]
 [132  94]]
```

	precision	recall	f1-score	support
0.0	0.73	0.93	0.82	389
1.0	0.78	0.42	0.54	226
micro avg	0.74	0.74	0.74	615
macro avg	0.76	0.67	0.68	615
weighted avg	0.75	0.74	0.72	615

Création d'un pipeline complet pour sauvegarder le modèle et le tester sur de nouvelles données.

In [77]:

```
1  from sklearn.pipeline import Pipeline
2
3  ▼ names=[
4      "NumTimesPrg", "PlGlcConc", "BloodP",
5      "SkinThick", "TwoHourSerIns", "BMI",
6      "DiPedFunc", "Age", "HasDiabetes"]
7  #il faut sauter les 9 premières lignes qui sont le descriptif des variables
8  df = pd.read_csv('pima-indians-diabetes.csv', names=names, skiprows=9)
9
10 array = df.values
11 X = array[:,0:8]
12 y = array[:,8]
13
14
15 print ('Création du pipeline \n')
16 ▼ pipeline = Pipeline([('scl', StandardScaler()),
17                        ('clf', svc)])
18
19 validation_size=0.3 #30% du jeu de données pour le test
20
21 testsize= 1-validation_size
22 seed=30
23 ▼ X_train,X_test,y_train,y_test=train_test_split(X, y,
24                                                  train_size=validation_size,
25                                                  random_state=seed,
26                                                  test_size=testsize)
27
28
29 pipeline.fit(X_train, y_train)
30 result = pipeline.predict(X_test)
31
32 print('\n accuracy:', accuracy_score(result, y_test), '\n')
33
34 import pickle
35 filename = 'modelsvcpima.pkl'
36 pickle.dump(pipeline, open(filename, 'wb'))
```

Création du pipeline

accuracy: 0.7881040892193308

In [78]:

```
1 print ("Chargement du modèle \n")
2 filename = 'modelsvcpima.pkl'
3 clf_loaded = pickle.load(open(filename, 'rb'))
4
5 #Considérons deux nouvelles données qui ne sont pas standardisées
6 ▼ new_df = pd.DataFrame([[6, 168, 72, 35, 0, 43.6, 0.627, 65],
7                           [2,85,67,30,2,27.6,0.351,22]], dtype='float')
8 ▼ print ("La première ligne correspond à une personne qui a tendance à avoir
9         "la prédiction devrait être de 1. Pour la seconde elle devrait être
10        "il faut remarquer que les données ne sont pas standardisées à l'entr
11        "de la prédiction une étape de standardisation à lieu grâce au pipel
12 prediction = clf_loaded.predict(new_df)
13 print(prediction)
14
```

Chargement du modèle

La première ligne correspond à une personne qui a tendance à avoir du diabète de type 2 la prédiction devrait être de 1. Pour la seconde elle devrait être de 0 il faut remarquer que les données ne sont pas standardisées à l'entrée et que lors de la prédiction une étape de standardisation à lieu grâce au pipeline

[1. 0.]

Essai de df2 sans standardisation.

Nous avons vu que SVC avait de meilleurs résultats aussi nous l'utilisons ici pour voir l'intérêt de la standardisation.

In [67]:

```
1 array = df2.values
2 X = array[:,0:8]
3 y = array[:,8]
4
5 svc = gd_sr.best_estimator_
6
7 validation_size=0.2 #20% du jeu de données pour le test
8
9 testsize= 1-validation_size
10 seed=20
11 X_train,X_test,y_train,y_test=train_test_split(X, y,
12                                                train_size=validation_size,
13                                                random_state=seed,
14                                                test_size=testsize)
15
16
17 svc.fit(X_train, y_train)
18 result = svc.predict(X_test)
19 print('\n accuracy: ', accuracy_score(result, y_test),'\n')
20
21 conf = confusion_matrix(y_test, result)
22 print ('\n matrice de confusion \n',conf)
23 print ('\n',classification_report(y_test, result))
```

accuracy: 0.7382113821138211

matrice de confusion

```
[[366  23]
 [138  88]]
```

	precision	recall	f1-score	support
0.0	0.73	0.94	0.82	389
1.0	0.79	0.39	0.52	226
micro avg	0.74	0.74	0.74	615
macro avg	0.76	0.67	0.67	615
weighted avg	0.75	0.74	0.71	615

Comme attendu, étant donné qu'il y a des données à des échelles très différentes l'accuracy diminue lorsque les données ne sont pas transformées.