


```

scale_Time = scaler.fit_transform(df[["scale_Time"]])
fl_liste1 = [item %>% subselect in scale_Time.collist() %>% item in subselect]
scale_Time = pd.Series(fl_liste1)
scale_Time

0    -1.997
1    -1.997
2    -1.997
3    -1.997
4    -1.997
...
284802    1.642
284803    1.642
284804    1.642
284805    1.642
284806    1.642
284807    1.642
Length: 284807, dtype: float64

# Mise à l'échelle de la colonne "Amount"
scale_Amount = scaler2.fit_transform(df[["Amount"]])
fl_liste2 = [item %>% subselect in scale_Amount.collist() %>% item in subselect]
scale_Amount = pd.Series(fl_liste2)
scale_Amount

0     0.245
1    -0.342
2     1.141
3     1.141
4    -0.370
...
284802    -0.370
284803    -0.370
284804    -0.370
284805    -0.370
284806    -0.370
284807    -0.370
Length: 284807, dtype: float64

On peut maintenant concatener les colonnes nouvellement créées dans le data-frame de départ le df

df = pd.concat([df, scale_Amount.rename("scale_Amount"), scale_Time.rename("scale_Time")],
axis=1)
df.sample(10)

V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V22      V23      V24      V25      V26      V27      V28      Class
16066  74936.0  1.335  0.331  -2.058  -0.346  2.583  2.854  -0.188  0.685  ...  -0.287  ...  -0.114  0.916  0.730  0.384  -0.032  0.030  0.76
17342  121447.0  1.986  -1.247  0.164  -0.208  -1.745  -0.064  -1.748  0.180  ...  1.070  ...  -0.101  -0.039  -0.325  -0.029  0.094  0.002  46.99
5513   47364.0  1.343  0.036  -0.376  -0.167  0.270  -0.114  0.100  -0.095  ...  -0.020  ...  -0.248  -0.730  0.692  1.144  -0.095  -0.028  2.00
95789  65467.0  0.883  -1.026  1.129  -0.800  -0.103  0.440  -1.344  0.277  ...  0.952  ...  -0.341  -0.077  0.402  -0.041  0.029  0.060  243.00
15635  107812.0  0.217  1.331  -0.834  0.186  0.419  -1.410  0.502  -0.006  ...  1.045  ...  -0.105  -0.068  -0.278  -0.193  -0.089  -0.052  15.0
19425  65092.0  -0.326  0.692  1.422  -0.681  0.157  -0.964  1.205  -0.406  ...  -0.219  ...  -0.081  0.510  -0.425  0.507  -0.293  -0.397  69.98
99780  131314.0  1.693  -0.706  -0.476  -0.425  -0.719  -0.531  -0.224  -0.131  ...  0.839  ...  -0.377  1.165  -0.704  0.154  -0.067  -0.004  161.93
22674  148842.0  0.514  1.380  1.091  1.605  0.461  -0.382  -0.599  -0.722  ...  0.567  ...  -0.130  0.769  0.910  -0.692  -0.033  0.097  16.70
17355  145206.0  2.170  -0.338  -2.580  0.766  0.665  -1.009  0.510  -0.373  ...  0.820  ...  -0.083  0.102  0.367  1.090  -0.161  -0.082  27.83
11109  79480.0  1.061  0.079  0.202  1.326  -0.249  -0.593  0.296  -0.118  ...  0.143  ...  -0.128  0.423  0.695  -0.273  0.014  0.023  63.96

10 rows x 31 columns

Ainsi on peut supprimer les anciennes colonnes "Amount" et "Time" qui seront remplacées par les colonnes
nouvellement créées ("scale_Amount" et "scale_Time")

df.drop(["Amount", "Time"], axis=1, inplace=True)
df.sample(1)

V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V22      V23      V24      V25      V26      V27      V28      Class
28154  1.021  -0.452  -2.203  0.879  -1.104  -0.223  1.198  -0.341  0.807  ...  -0.485  ...  -0.271  -0.154  -0.780  0.750  -0.287  0.035  0.005  0.0
16095  2.092  0.505  -2.708  0.800  0.862  -1.310  0.294  -0.294  0.674  ...  -1.184  ...  -0.222  1.021  0.028  -0.233  -0.236  -0.165  -0.030  0.0
109954  1.248  0.324  -0.968  1.320  2.385  3.708  -0.438  0.865  -0.866  ...  0.780  ...  -0.260  -0.106  1.001  0.729  0.083  -0.009  0.014  0.0
197781  1.332  0.328  -0.941  0.402  -0.648  -1.219  0.448  0.519  -0.320  ...  -1.387  ...  0.402  -0.049  -0.372  0.391  0.200  0.016  -0.015  0.0
234426  -1.066  3.084  -0.166  -0.087  -0.571  3.084  -3.952  -0.208  1.418  ...  -0.380  ...  -1.958  1.288  -1.205  -0.215  -0.015  0.593  0.166  0.0

5 rows x 31 columns

A priori, on peut diviser le dataset en deux : le train et le test. Le fichier d'entraînement doit être assez volumineux
pour obtenir des résultats statistiques significatifs. On choisit alors de sélectionner un échantillon aléatoire à
hauteur de 80% des données du fichier de base. Si le fichier est totalement aléatoire on risque d'avoir un fichier
d'entraînement contenant peu de transaction frauduleuse car ces dernières ne constituent que 0.17% du dataset df,
par conséquent notre modèle risque de ne pas être significatif. On décide alors de créer un nouveau fichier df2 qui
servira de fichier de base d'échantillonnage contenant 100% des transactions frauduleuses de df soit 50% de ces
données sont frauduleuses et 50% de celles qui ne sont pas frauduleuses (voir histogramme ci-dessous).

Ainsi on a d'autant de chance de sélectionner aléatoirement une transaction frauduleuse qu'une transaction. Ensuite
à partir de df2 on effectue un échantillonnage aléatoire pour obtenir un fichier d'entraînement et un fichier test.

# On commence d'abord créer le train et le test manuellement
DivCol = np.random.rand(n(df)) < 0.8 # un vecteur de booléen de taille len(df)
train = df[DivCol]
test = df[~DivCol]
train (forme du train: ()) Uniforme du test: ().format(train.shape, test.shape)

forme du train: (236162, 31)
forme du test: (48645, 31)

train.reset_index(drop=True, inplace=True) # faire que les indices commencent de 0 à len(train)
test.reset_index(drop=True, inplace=True) # faire que les indices commencent de 0 à len(test)

test

V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V22      V23      V24      V25      V26      V27      V28      Class
0      1.230  0.141  0.045  1.203  0.192  0.273  -0.005  0.841  0.607  ...  -0.485  ...  -0.271  -0.154  -0.780  0.750  -0.287  0.035  0.005  0.0
1      -0.782  -0.588  -0.908  -0.419  0.600  0.862  -1.310  0.294  -0.294  ...  0.674  ...  -0.222  1.021  0.028  -0.233  -0.236  -0.165  -0.030  0.0
2      0.862  0.328  -0.171  2.109  1.130  1.696  0.108  -0.522  -1.191  ...  0.724  ...  -0.402  -0.049  -0.372  0.391  0.200  0.016  -0.015  0.0
3      0.247  0.728  1.185  -0.003  -1.314  -0.150  -0.946  -1.618  1.544  ...  -0.830  ...  0.200  -0.185  0.423  0.821  -0.228  0.337  0.250  0.0
4      -1.452  1.765  0.612  1.177  -0.446  0.247  -0.258  1.092  -0.608  ...  0.047  ...  -0.328  -0.069  0.021  -0.045  -0.243  0.149  0.121  0.0
...
28740  -0.828  0.650  1.844  0.191  -0.605  0.580  -0.614  0.922  0.259  ...  -0.725  ...  1.052  -0.167  -0.445  -0.408  -0.065  0.104  0.017  0.0
28741  -0.765  0.588  -0.908  -0.419  0.600  0.862  -1.310  0.294  -0.294  ...  0.674  ...  -0.222  1.021  0.028  -0.233  -0.236  -0.165  -0.030  0.0
28
```


On constate qu'on a seulement 6 et 8 éléments qui sont mal étiquetés par le classificateur contre 123 et 54 éléments diagonaux qui sont bien étiquetés ce qui est suffisamment grand par-rapport à nos données du sous-ensemble. donc on peut dire que notre paramètre de régularisation est bon et peut-être meilleure.

```
# Import Modules
from sklearn.linear_model import LogisticRegression
from sklearn.decomposition import PCA
from sklearn.model_selection import GridSearchCV

# Create a default model
pipe_lr = Pipeline([('pca', PCA(n_component=1)),
                    ('l1r', LogisticRegression(penalty='l1',
                                                random_state=1,
                                                solver='lbfgs'))])

# Create a list of choices for each hyperparameter
param_range_pca = [2, 5, 10, 50] # dimension of the data after dimensionality reduction
param_clf = ['l1', 'l2']

# Create the grid of possible combinations
param_grid = [{'pca_n_components': param_range_pca,
               'clf_penalty': param_clf}]

# Prepare the Grid Search Fine Tuning
gs = GridSearchCV(estimator=pipe_lr, # the model
                  param_grid=param_grid, # its hyperparameters possibilities
                  scoring='precision', # precision is the evaluation metric
                  cv=5, # We use 5 folds in cross validation
                  n_jobs=-1)

# Test all the possible hyperparameters
gs = gs.fit(X_train, y_train)

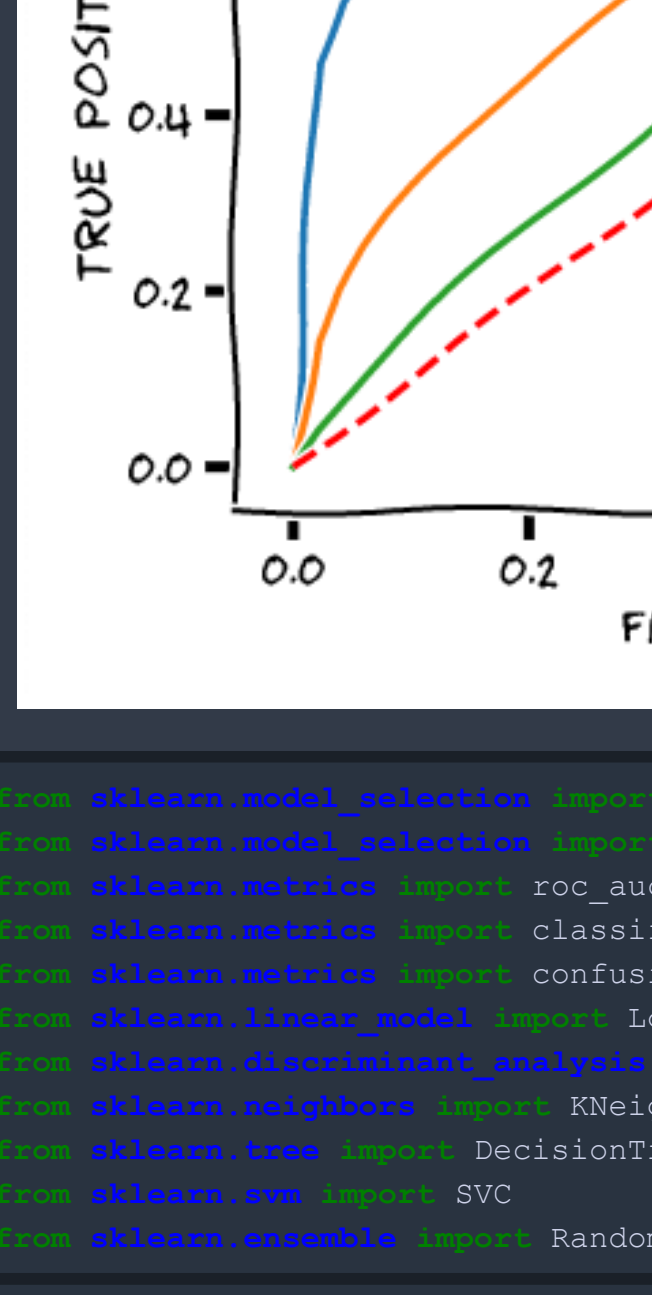
# The best hyperparameters
print("The optimal precision is {}".format(gs.best_score_))
print("The optimal hyperparameters are {}".format(gs.best_params_))

The optimal precision is 0.8614039796933468
The optimal hyperparameters are {'clf_penalty': 'l2', 'pca_n_components': 10}
```

```
# Optimal model
best_lr_svd = gs.best_estimator_

# Evaluation performances
crossvalidate(best_lr_svd, X_train, y_train)
# Confusion Matrix
plot_confusion_matrix(best_lr_svd, X_test, y_test)

Cross Validation Precision: 0.86 (+/- 0.03)
Cross Validation Recall: 0.86 (+/- 0.08)
Cross Validation F1 score: 0.91 (+/- 0.04)
Cross Validation roc_auc: 0.96 (+/- 0.03)
Cross Validation Accuracy: 0.94 (+/- 0.03)
```

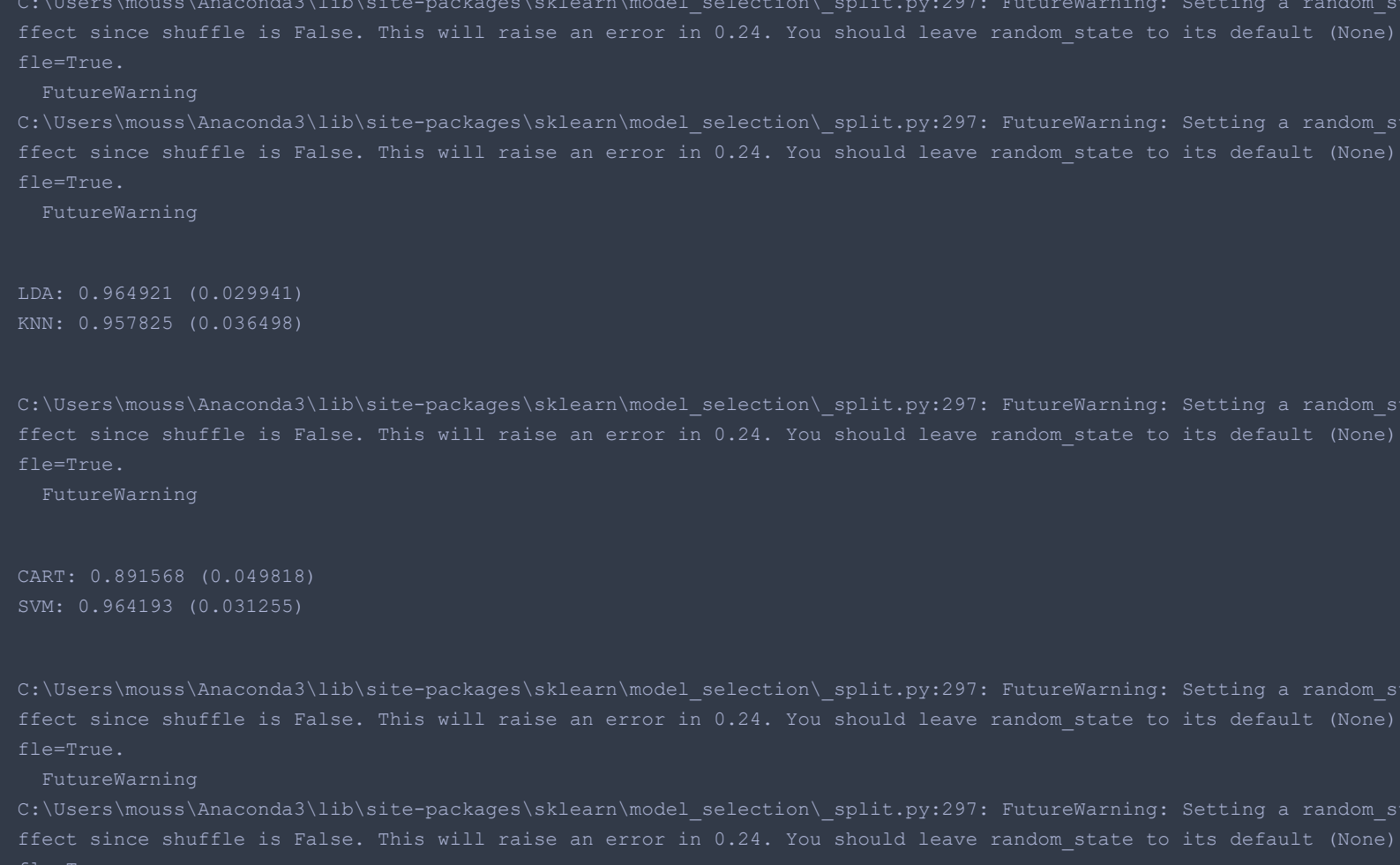


Ici, on constate que les éléments hors diagonaux sont réduits par-rapport à la matrice de confusion précédente. Donc le paramètre de régularisation est mieux qu'avant. Ainsi notre classifieur est meilleure que celui d'avant.

La courbe ROC est une courbe qui nous permettra de déterminer le seuil de décision adéquat pour notre modèle. Chaque point de la courbe représente le ratio sensibilité / spécificité pour un seul de décision spécifique. L'aire en dessous de la courbe ROC et nommé AUC et est comprise entre 0 et 1. L'AUC donne une mesure de la capacité du modèle à discriminer les cas positifs des cas négatifs.

Le schema ci dessous nous montre quel est la forme idéal et standard de la courbe.

```
from IPython.display import Image
Image(filename = "image1.png", width=1200, height=1200)
```



```
from sklearn.model_selection import KFold
from sklearn.metrics import cross_val_score
from sklearn.metrics import roc_auc_score
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.linear_model import LogisticRegression
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier

#On peut créer les modèles par défauts
modeles = []
modeles.append(('LDA', LinearDiscriminantAnalysis()))
modeles.append(('KNN', KNeighborsClassifier()))
modeles.append(('CART', DecisionTreeClassifier()))
modeles.append(('SVM', SVC()))
modeles.append(('RF', RandomForestClassifier()))
modeles.append(('LR', LogisticRegression()))
```

```
# On peut tester les modèles
resultats = []
list_modeles = []
for i, modele in modeles:
    kfold = KFold(n_splits=10, random_state=40)
    cv_resultats = cross_val_score(modele, X_train, y_train, cv=kfold, scoring='roc_auc')
    resultats.append(cv_resultats)
    list_modeles.append(i)
    message = '%s %f %f' % (i, cv_resultats.mean(), cv_resultats.std())
    print(message)

C:\Users\mouse\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:297: FutureWarning: Setting a random_state has no e
fect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuf
file=True.
FutureWarning
C:\Users\mouse\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:297: FutureWarning: Setting a random_state has no e
fect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuf
file=True.
FutureWarning

IDN: 0.964921 (0.029941)
FNN: 0.957825 (0.036698)

C:\Users\mouse\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:297: FutureWarning: Setting a random_state has no e
fect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuf
file=True.
FutureWarning

CART: 0.891368 (0.048818)
SVM: 0.864193 (0.032250)

C:\Users\mouse\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:297: FutureWarning: Setting a random_state has no e
fect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuf
file=True.
FutureWarning
FutureWarning

RF: 0.971546 (0.027922)

C:\Users\mouse\Anaconda3\lib\site-packages\sklearn\model_selection\_split.py:297: FutureWarning: Setting a random_state has no e
fect since shuffle is False. This will raise an error in 0.24. You should leave random_state to its default (None), or set shuf
file=True.
FutureWarning

LR: 0.963628 (0.029673)
```

On ainsi comparer les algorithmes de classifications en faisant un boxplot des modèles

```
fig = plt.subplots(figsize = (10,10))
plt.title("Comparaison des différents algorithmes de classifications", size=20)
ax.set_xticklabels(list_modeles)
plt.boxplot(resultats)
plt.xlabel("Algorithme", size=10)
plt.ylabel("ROC-AUC Score", size=20)

C:\Users\mouse\Anaconda3\lib\site-packages\ipykernel_launcher.py:3: UserWarning: FixedFormatter should only be used together wit
h FixedLocator
This is separate from the ipykernel package so we can avoid doing imports until

Test(0, 0.5, 'ROC-AUC Score')
```



On remarque le meilleur classifieur est RandomForestClassifier.