

2023

# **RAPPORT OPTIMISATION D'UN MODÈLE**

**Encadrée par : Simplon , Sanofi**

**Lecarpentier Eden**



# SOMMAIRE

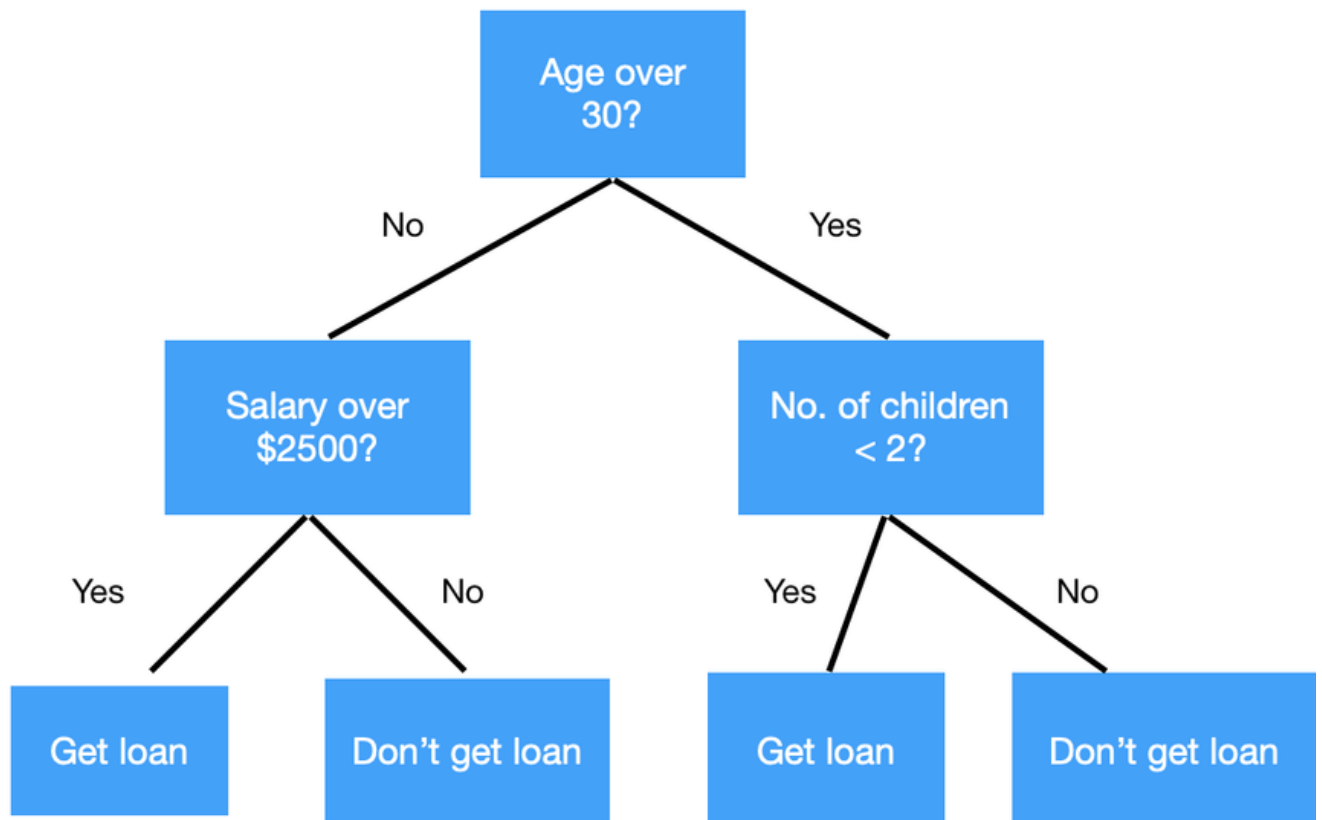
- 01**    Modèle de référence
- 02**    Performance de référence
- 03**    Optimisation du modèle
- 04**    GridSearchCV
- 05**    SMOTE
- 06**    modèle optimisé performance
- 07**    Test unitaire sur l'application
- 08**    Conclusion

# MODÈLE DE RÉFÉRENCE

Le projet que je vais améliorer est sur la prédiction de problème cardiaques .  
On a une target contenant 2 classes Oui et Non .  
C'est un modèle de classification en machine learning .  
Mon modèle a été réalisé avec un RandomForest classifieur .

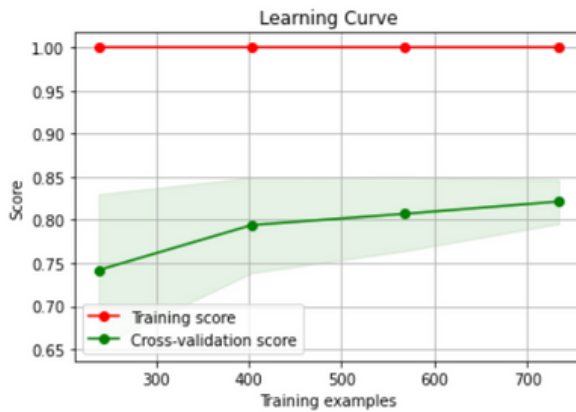
## RandomForestClassifier

- Model de classification
- Le random forest fonctionne comme un arbre .
- On part du tron et on va dans les branches
- Cela permet de trouver toutes les corrélations entre les features

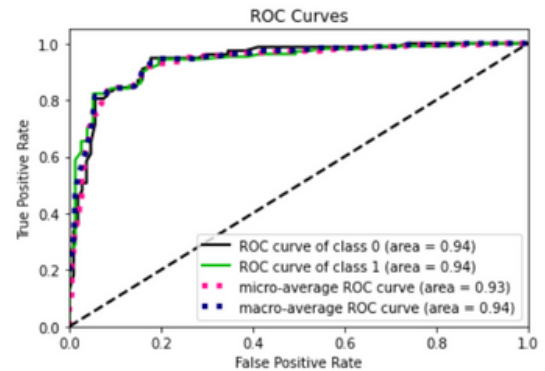


# MODÈLE DE RÉFÉRENCE

## LearningCurve



## ROC CURVE



- La Learning Curve permet de voir les compétences du modèle sur le train (ligne rouge) et le test (ligne verte)
- Elle permet de trouver de l'overfitting ou Underfitting, le taux d'erreur possible
- L'espace des deux courbes montre le taux d'erreur possible.
- Les courbes sont probablement trop éloignées pour que le modèle soit efficace
- Les courbes risquent également de se toucher

- La ROC Curve permet de vérifier l'entraînement du modèle, plus les lignes montent vite mieux sera l'entraînement.
- Les courbes représentent 2 paramètres : le taux de vrai et faux positif.
- Cette ROC Curve montre une bonne performance

## Classification Report

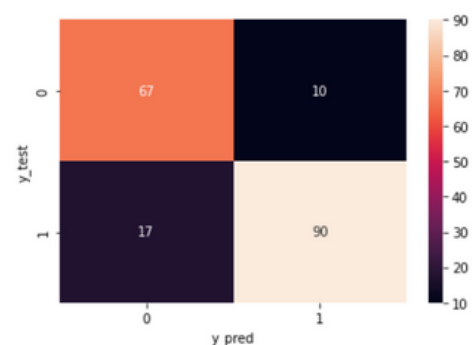
accuracy train : 1.000

accuracy test : 0.891

	precision	recall	f1-score	support
0	0.76	0.88	0.82	77
1	0.91	0.80	0.85	107
accuracy			0.84	184
macro avg	0.83	0.84	0.84	184
weighted avg	0.85	0.84	0.84	184

- La classification report nous permet de voir les performances sur les classes que l'on veut prédire grâce au metric.
- Nos classes : oui cette personne est à risque, non cette personne n'est pas à risque d'une crise cardiaque
- Le modèle de classification utilise les metrics : accuracy et recall
- L'accuracy n'est pas mauvaise par contre le recall la classe 0 (Non) devrait se rapprocher le plus possible de 0 et l'inverse pour la classe 1 (Oui)

## Confusion Matrix



- La Confusion Matrix nous montre : les faux positifs, les faux négatifs, les vrais positifs, les vrais négatifs
- Vrai positif : devait avoir une crise cardiaque et l'a eue
- Faux positif : devait ne pas avoir de crise mais en a eue une
- Vrai négatif : devait ne pas avoir de crise et n'en a eue aucune
- Faux négatif : devait ne pas avoir de crise et en a eue une

# OPTIMISATION DU MODÈLE

L'optimisation d'un modèle se fait en plusieurs étapes pour ce projet j'ai réalisé les choses suivantes:  
Trouvée de bons hyperparamètres, équilibrés les données.

## GridSearchCV

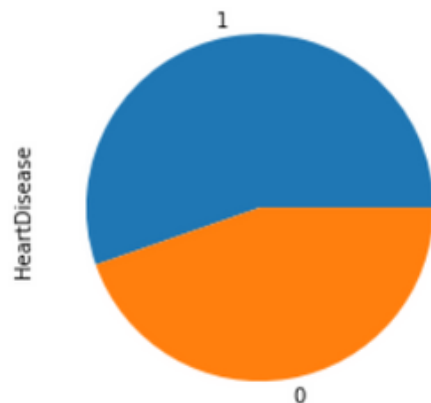
- GridSearchCV est utilisée pour trouver les meilleurs hyperparamètres du modèle à utiliser pour normalement améliorer les performances du modèle basé sur les données utilisées.
- Les hyperparamètres que je vais utiliser pour la régression logistique sont les suivants : C, penalty, solver.
- Les valeurs conseillées par le GridSearchCV sont : n\_estimators, max\_features, max\_depth, criterion.
- Hyperparamètre : ce sont des paramètres qui permettent de déterminer les valeurs nécessaires pour le paramètre du modèle.
- n\_estimators : Indication sur le nombre d'arbres dans la forêt j'en utilise 200
- max\_features : Détermine le nombre maximal de features possible pour chaque arbre
- max\_depth : Le nombre de split possible pour chaque arbre permet de lutter contre l'overfit, underfit
- criterion : Permet de vérifier la qualité d'un split en utilisant Gini Impurity qui permet de savoir à quel fréquence des éléments du set sont indiqués incorrectement.

## SMOTE

Smote est utilisée lorsque la cible n'est pas équilibrée, pour déterminer si on utilise le SMOTE ou non ou utiliser un graphique :

Le graphique suivant nous montre que les 2 classes ne sont pas égales.

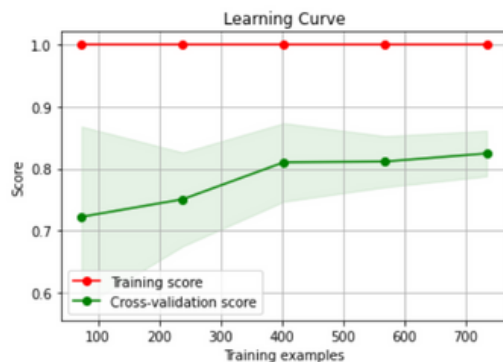
Les 2 classes étant si oui ou non on aura une crise cardiaque.



L'utilisation du SMOTE permettra de créer des données synthétiques c'est-à-dire des données qui diffèrent légèrement des données originales. Cela permet d'avoir plus de données d'une des classes qui en manque donc équilibrer les données pour essayer d'améliorer le modèle.

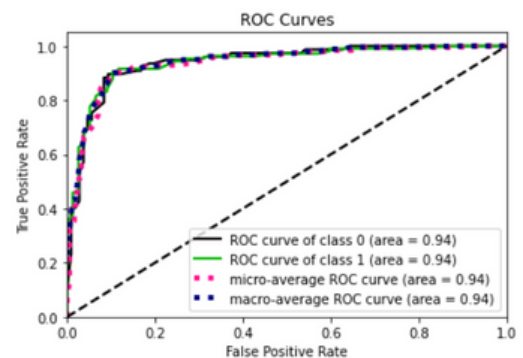
# MODÈLE OPTIMISÉ PERFORMANCE

## LearningCurve



- La Learning Curve est toujours en overfitting
- La courbe test s'entraîne légèrement mieux

## ROC CURVE



- La Roc Curve est resté relativement similaire a la précédente .
- On peut s'apercevoir que les lignes montes de façon régulière contrairement à la précédente Roc Curve .

## Classification Report

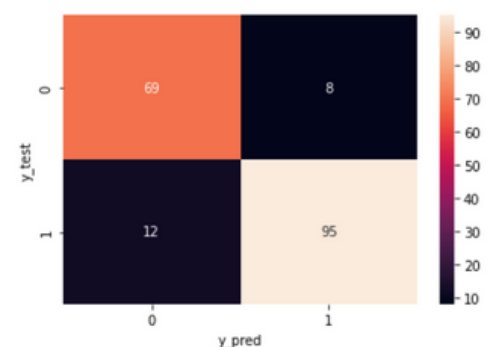
accuracy train : 0.857

accuracy test : 0.837

	precision	recall	f1-score	support
0	0.80	0.87	0.83	77
1	0.90	0.84	0.87	107
accuracy			0.85	184
macro avg	0.85	0.86	0.85	184
weighted avg	0.86	0.85	0.85	184

La class 1 a augmenter de 0.04 passant de 0.80 à 0.84.  
La classe 0 a baisser de 0.01 passant de 0.88 à 0.87.

## Confusion Matrix



La confusion matrix a plus de True positive est le reste a diminué .

# TEST UNITAIRE SUR L'APPLICATION

## Application

The first screenshot shows the 'User Input Features' form with fields for age, sex, chest pain type, resting blood pressure, serum cholesterol, fasting blood sugar, resting electrocardiographic results, maximum heart rate achieved, exercise-induced angina, oldpeak, and the slope of the peak exercise ST segment. The second screenshot shows the 'Heart disease Prediction App' interface with a table of 10 input rows and a 'Prediction' button. The third screenshot shows the 'Prediction' result table with 10 rows of output values.

- Ajout d'un bouton
- Le bouton permet de montr  les pr diction une fois le bouton cliqu  est non avant

## Test unitaire

1. V rification que mon dataset est le bon nombres de colonnes .
2. V rification que les pr diction de l'application reste entre 0 et 1 si mon mod le aller au dessus ou en dessous de ces nombres cela voudrais dire que mon mod le est d fectueux .
3. V rification que le nombre de pr diction sois  gale au nombre d' chantillons donner

test\_main.py ...

[100%]

===== 3 passed in 1.56s =====

# CONCLUSION

Le projet a eu du mal à être optimisée probablement du à un grand manque de donnée car la base de données est petite.

Le modèle pourrait être encore améliorée en ajoutant de nouvelles données et de la diversité dans les données.

Cette diversité des données est nécessaire pour que le modèle soit le plus performant possible sur des données qu'il n'a jamais vues.

Le problème le plus voyant est l'overfitting même après l'avoir optimisée la conclusion amenée précédemment aiderait le modèle à s'améliorer d'avantages.