

IRANIAN CHURN PREDICTION

PYTHON FOR DATA ANALYSIS – A5

CAZIER SIMON - CHASLE JEAN-LUC

SOMMAIRE

- PRESENTATION DU PROBLÈME
- PROBLÉMATIQUE
- COMPRÉHENSION ET NETTOYAGE DES DONNÉES
- VISUALISATION DES DONNÉES
- MODÉLISATION DES DONNÉES
- API FLASK



PRÉSENTATION DU PROBLÈME

- FACE À LA CONCURRENCE RUDE DANS UN DOMAINE EN PERPÉTUEL EXPANSION, CELUI DE LA TÉLÉCOMMUNICATION, LES OPÉRATEURS DOIVENT SAVOIR RECONNAÎTRE ET ANALYSER LES ATTITUDES DE LEURS CLIENTS.
- POUR ÊTRE COMPÉTITIVES, LES ORGANISATIONS DOIVENT CONNAÎTRE ET PRÉVOIR LES PRÉFÉRENCES ET LES COMPORTEMENTS DES CLIENTS AFIN DE MAXIMISER LE NOMBRE DE FIDÉLISATIONS AVANT QUE LEURS CONCURRENTS NE LE FASSENT OU QUE CEUX-CI RÉSILIENT VERS UN AUTRE OPÉRATEUR.
- DANS CETTE ÉTUDE, ON TENTERA D'IDENTIFIER LES FACTEURS (PARAMÈTRES) QUI INFLUENT SUR LE TAUX DE RÉSILIATION (COLONNE CHURN) AFIN DE PERMETTRE À L'OPÉRATEUR DE CONSERVER SES CLIENTS ET D'EN ACQUÉRIR DE NOUVEUX, LE PLUS PRÉCIEUX DES ACTIFS D'UNE ORGANISATION .

PROBLÉMATIQUE

LE BUT DE NOTRE PROJET EST DE PRÉDIRE LA RÉSILIATION POTENTIELLE D'UN CLIENT OU/ET D'Étudier LES PARAMÈTRES QUI CARACTÉRISent LA RÉSILIATION DES CLIENTS.

COMPRÉHENSION ET NETTOYAGE DES DONNÉES

- LES DONNÉES REGROUPENT 3150 CLIENTS SÉLECTIONNÉS DE FAÇON ALÉATOIRE DANS LA BASE DE DONNÉES D'UNE COMPAGNIE DE TÉLÉCOMMUNICATION IRANIENNE SUR UN AN.

- 14 VARIABLES DÉCRIVENT LE COMPORTEMENT ET LE STATUT DU CLIENT AVEC UN LABEL CHURN QUI SERA À PREDIRE SI UN CLIENT VA RÉSILIER OU NON

Nom de la variable	Description
Complains	Si le client s'est plaint au moins une fois auprès du réseau de télécommunications
Subscription Length	Le nombre total de mois d'abonnement du client
Charge Amount	Prix de l'abonnement allant du moins cher au plus cher sous forme de catégories
Seconds of use	Le nombre total de secondes de tous les appels d'un client
Frequency of use	Le nombre total d'appels d'un client
Frequency of SMS	Le nombre total de sms d'un client
Distinct Called Numbers	Le nombre total d'appels téléphoniques distincts d'un client
Age Group	Catégorie du client selon son âge
Tariff Plan	Si le forfait est contractuel ou si c'est une carte prépayée
Status	Si le client est actif ou inactif
Customer Value	La valeur du client auprès de l'entreprise (note)
Age	Age du clie
Call Failures	Nombre d'appels échoués
Churn	Variable à prédire (1 = résilier, 0 = Toujours abonné)

COMPRÉHENSION ET NETTOYAGE DES DONNÉES

- APRÈS AVOIR EXPLORER LES DONNÉES, ON CONSTATE QU'IL N'Y A AUCUNE VALEUR MANQUANTE

Vérification s'il existe des valeurs NULL présents dans le dataset

```
nb_null_colonne = df_iranian_churn.isnull().sum()
```

```
nb_null_colonne
```

Call Failure	0
Complains	0
Subscription Length	0
Charge Amount	0
Seconds of Use	0
Frequency of use	0
Frequency of SMS	0
Distinct Called Numbers	0
Age Group	0
Tariff Plan	0
Status	0
Age	0
Customer Value	0
Churn	0
dtype:	int64

Il n'y a pas de valeurs manquantes dans le dataset donc il n'y pas besoin de gérer les valeurs null

VISUALISATION DES DONNÉES

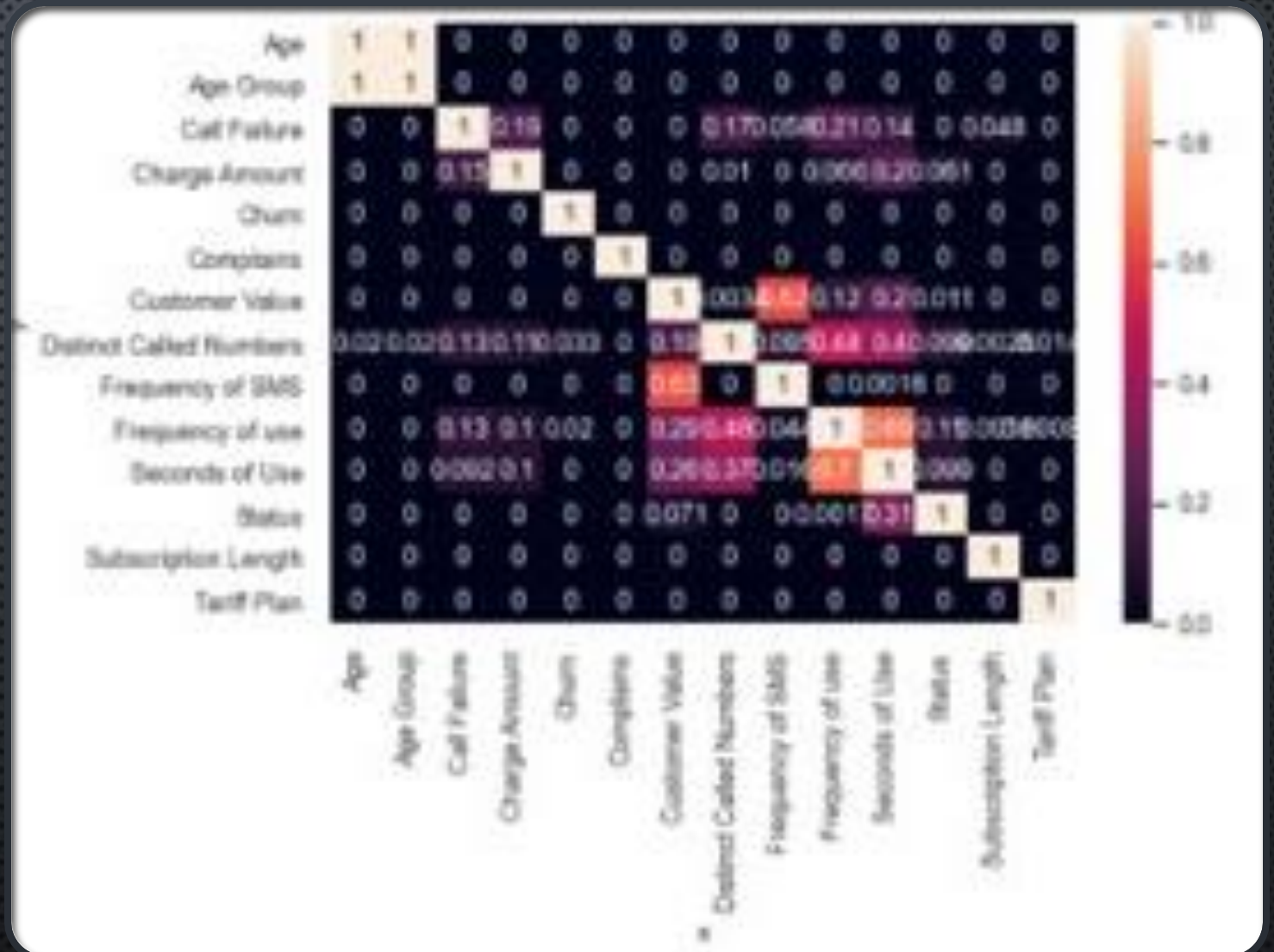
- LA MATRICE DE CORRÉLATION NOUS DONNE L'IMPACT DES DIFFÉRENTS PARAMÈTRES SUR LA VARIABLE "CHURN".

- ON PEUT CONSTATER QUE LES VARIABLES "COMPLAINS" ET "STATUS" SONT CELLES QUI INFLUENT LE PLUS SUR LE LABEL DE RÉSILIATION.

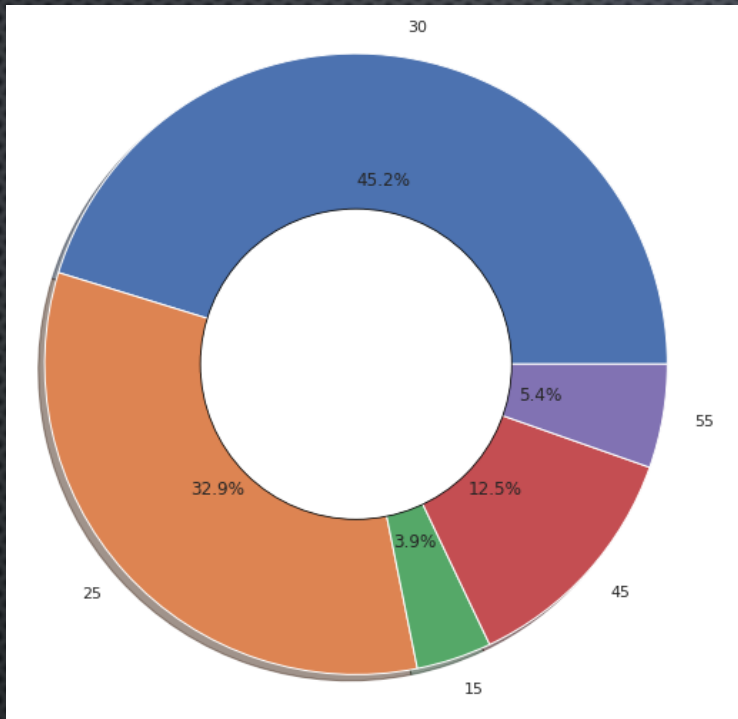


ALTERNATIVE À LA MATRICE DE CORRÉLATION : PPScore

- LE PPS EST UN SCORE ASYMÉTRIQUE, DE TYPE DONNÉES-DIAGNOSTIC, QUI PERMET DE DÉTECTER DES RELATIONS LINÉAIRES OU NON LINÉAIRES ENTRE DEUX COLONNES. LE SCORE VA DE 0 (AUCUN POUVOIR PRÉDICTIF) À 1 (POUVOIR PRÉDICTIF PARFAIT). IL PEUT ÊTRE UTILISÉ COMME ALTERNATIVE À LA CORRÉLATION (MATRICE).
- UN SCORE DE 0 SIGNIFIE QUE LA COLONNE X NE PEUT PAS PRÉDIRE LA COLONNE Y MIEUX QU'UN MODÈLE DE BASE NAÏF.
- UN SCORE DE 1 SIGNIFIE QUE LA COLONNE X PEUT PARFAITEMENT PRÉDIRE LA COLONNE Y ÉTANT DONNÉ LE MODÈLE.
- UN SCORE ENTRE 0 ET 1 INDIQUE LE RAPPORT ENTRE LA PUISSANCE PRÉDICTIVE POTENTIELLE DU MODÈLE ET CELLE DU MODÈLE DE BASE.

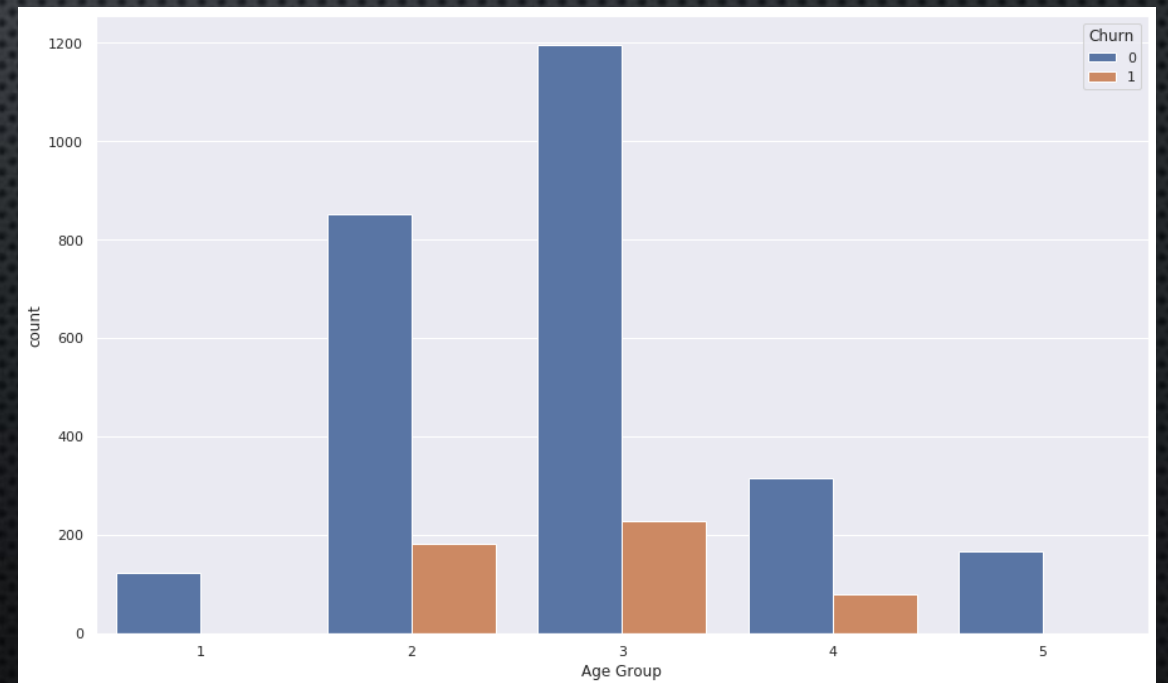


Distinct called number et frequency of use semblent potentiellement prédire la variable cible

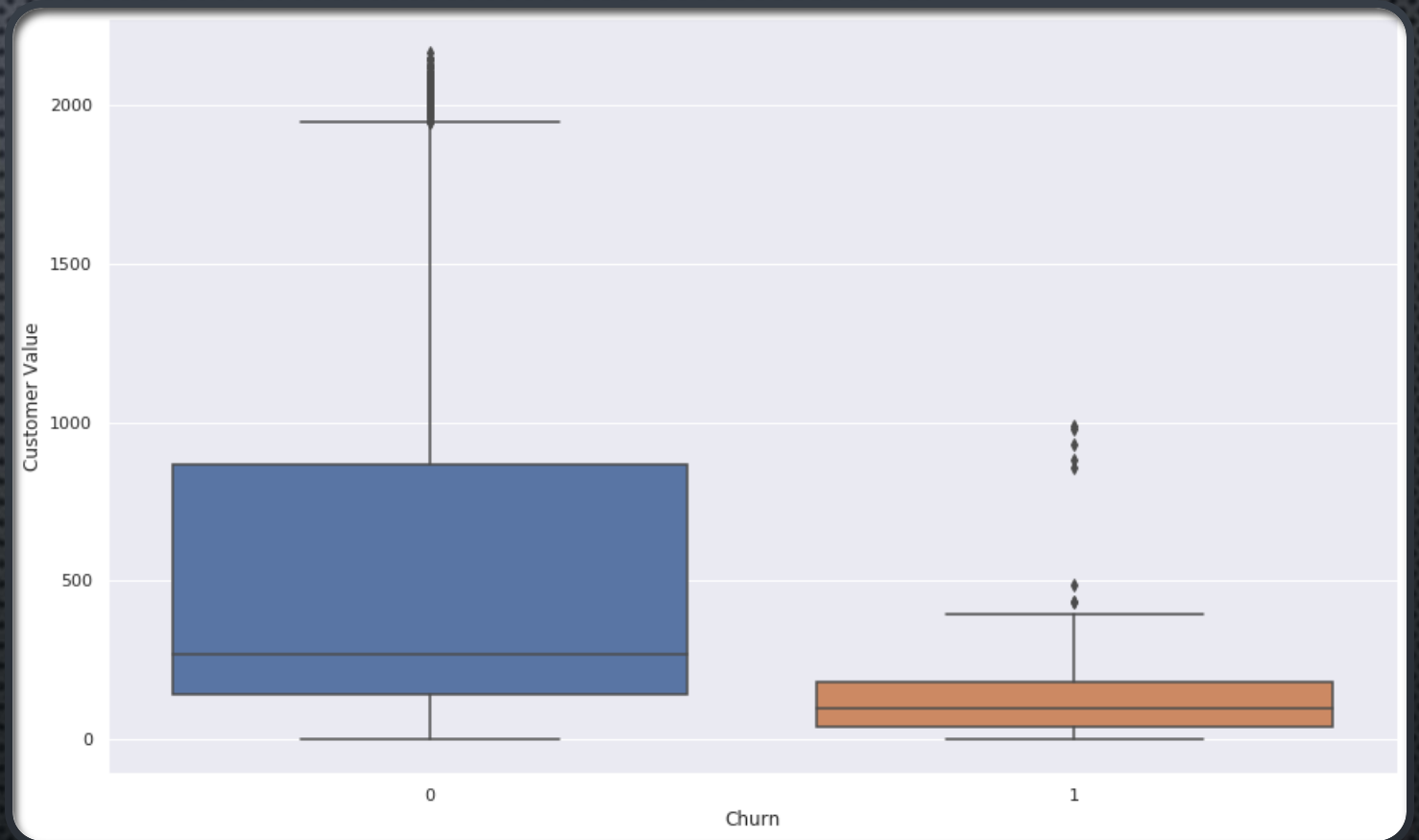


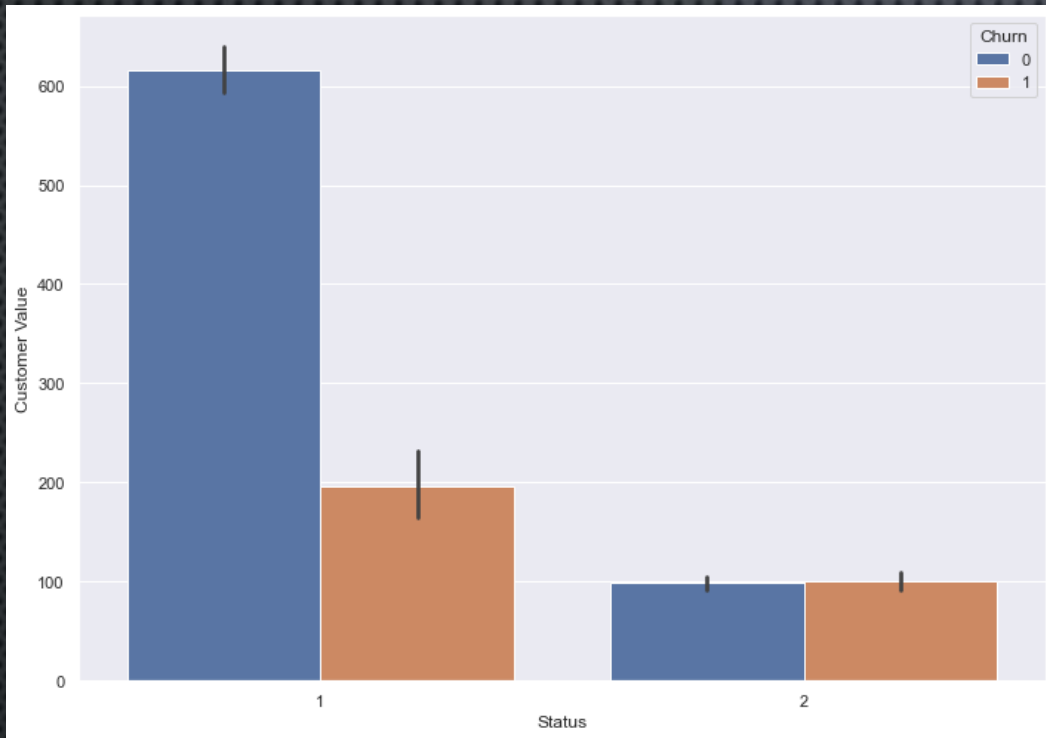
Ce second graphe représente le nombre de résiliations de clients par groupe d'âge. On peut constater que les clients entre 25 et 45 ont plus tendance à résilier leur abonnement que le reste des autres clients.

Ce graphe représente la répartition de la clientèle en fonction de leur âge au sein du dataset. La majorité des utilisateurs sont des personnes entre 25 et 30 ans.

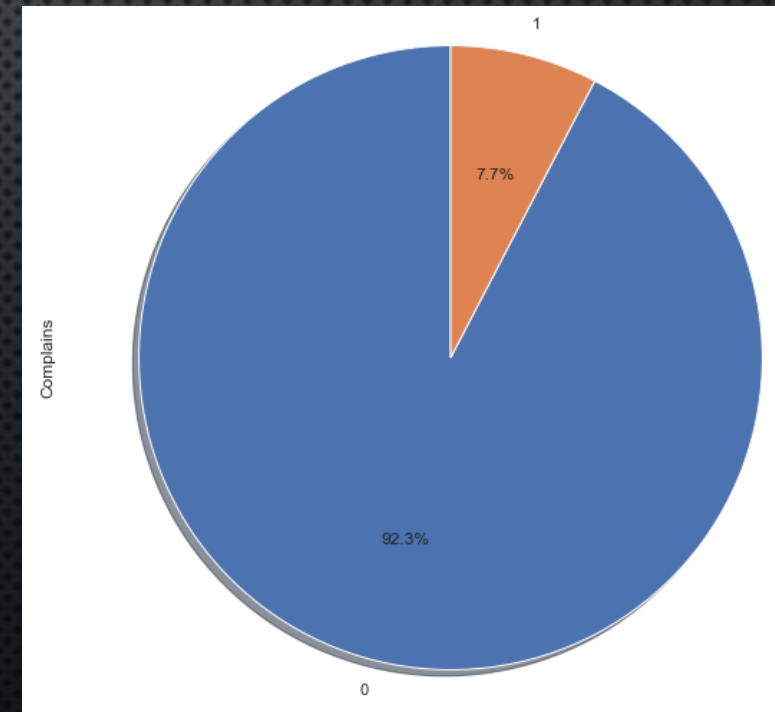


- Le graphique de droite nous indique que le "Customer Value" (valeur du client auprès de l'entreprise de télécommunication) a un impact sur la prise de décision d'un client lors de la résiliation de son abonnement.
- En dessous d'un seuil d'environ 200, la quasi-totalité des clients résilient.
- Ceci est donc un bon indicateur pour l'entreprise afin de connaître les clients qui vont résilier.





Ce graphe circulaire nous informe que 7.7% des clients qui se sont plaint vont jusqu'à résilier leur abonnement.



Cette représentation du "Customer Value" en fonction des deux paramètres "Status" et "Churn", nous indique que presque tous les utilisateurs non-actifs (status: 2) résilient leurs abonnements

MODÉLISATION DES DONNÉES : MACHINE LEARNING

- AFIN DE PRÉDIRE LA VALEUR CHURN D'UN CLIENT, NOUS ALLONS UTILISER 5 ALGORITHMES DE CLASSIFICATION PUISQUE LA VALEUR CIBLE EST CATÉGORIQUE (VALEURS DISCRÈTES).
- MODÈLES UTILISÉS :
 1. LOGISTIC REGRESSION
 2. LOGISTIC REGRESSION WITH SMOTE
 3. RANDOM FOREST
 4. KNN
 5. GRADIENT BOOSTING

ON A DÉCOUPÉ PAR LA SUITE DE NOTRE JEU DE DONNÉES EN UN JEU D'ENTRANEMENT (80% = 2520 LIGNES) ET UN JEU DE TEST (20% = 630 LIGNES)

LOGISTIC REGRESSION

Avec la régression logistique, nous avons une précision de 0.913.

Une précision de 0.92 et de 0.88 pour les valeurs 0 et 1 respectivement.

L'algorithme prédit de manière satisfaisante les labels mais on a un nombre assez conséquent de FP = 48, qui est un danger pour le besoin métier, en effet, on prédit qu'un client ne va pas résilier alors qu'il va résilier ce qui n'arrange pas l'opérateur

Accuracy Score:

0.913

Classification report:

	precision	recall	f1-score	support
0	0.92	0.99	0.95	532
1	0.88	0.51	0.65	98
accuracy			0.91	630
macro avg	0.90	0.75	0.80	630
weighted avg	0.91	0.91	0.90	630

Confusion matrix:

```
[[525  7]
 [ 48 50]]
```


LOGISTIC REGRESSION WITH SMOTE

- LA CLASSIFICATION DÉSÉQUILBRÉE IMPLIQUE LE DÉVELOPPEMENT DE MODÈLES PRÉDICTIFS SUR DES ENSEMBLES DE DONNÉES DE CLASSIFICATION QUI PRÉSENTENT UN GRAVE DÉSÉQUILIBRE DE CLASSE.
- LE DÉFI DE TRAVAILLER AVEC DES ENSEMBLES DE DONNÉES DÉSÉQUILIBRÉS EST QUE LA PLUPART DES TECHNIQUES D'APPRENTISSAGE AUTOMATIQUE IGNORENT LA CLASSE MINORITAIRE ET ONT À LEUR TOUR DE MAUVAISES PERFORMANCES SUR CETTE CLASSE, BIEN QUE CE SOIT GÉNÉRALEMENT LA PERFORMANCE SUR LA CLASSE MINORITAIRE QUI EST LA PLUS IMPORTANTE.
- UNE APPROCHE POUR REMÉDIER AUX ENSEMBLES DE DONNÉES DÉSÉQUILIBRÉS CONSISTE À SURÉCHANTILLONNER LA CLASSE MINORITAIRE. L'APPROCHE LA PLUS SIMPLE CONSISTE À DUPLIQUER DES EXEMPLES DANS LA CLASSE MINORITAIRE, BIEN QUE CES EXEMPLES N'AJOUTENT AUCUNE NOUVELLE INFORMATION AU MODÈLE. AU LIEU DE CELA, DE NOUVEAUX EXEMPLES PEUVENT ÊTRE SYNTHÉTISÉS À PARTIR DES EXEMPLES EXISTANTS. IL S'AGIT D'UN TYPE D'AUGMENTATION DES DONNÉES POUR LA CLASSE MINORITAIRE, APPELÉ "SYNTHETIC MINORITY OVERSAMPLING TECHNIQUE" (SMOTE).

LOGISTIC REGRESSION WITH SMOTE

- Avec le SMOTE, la regression logistique a une précision plus faible pour le nombre de clients ayant résiliés que celle sans. Le surechantillonnage du nombre d'échantillons pour la classe 1 (le client résilie) a entraîné plus de FP qui est moins grave que d'avoir plus de FN.
- Par ailleurs, on constate qu'on a moins de FN ce qui est encourageant et on prédit mieux les personnes qui souhaitent résiliées avec cette méthode.

Accuracy Score:

0.837

Classification report:

	precision	recall	f1-score	support
0	0.98	0.82	0.89	521
1	0.52	0.91	0.66	109
accuracy			0.84	630
macro avg	0.75	0.86	0.78	630
weighted avg	0.90	0.84	0.85	630

Confusion matrix:

[[428 93]

[10 99]]

RANDOM FOREST

Avec la régression logistique, nous avons une précision de 0.946.

Une précision de 0.96 et de 0.89 pour les valeurs 0 et 1 respectivement.

L'algorithme prédit de manière satisfaisante les labels et fait très peu d'erreur.

C'est un algorithme que l'on pourrait facilement retenir pour le besoin de l'opérateur

Accuracy Score:

0.948

ROC Score:

0.9828487911391292

Classification Report:

	precision	recall	f1-score	support
0	0.96	0.98	0.97	521
1	0.89	0.80	0.84	109
accuracy			0.95	630
macro avg	0.92	0.89	0.90	630
weighted avg	0.95	0.95	0.95	630

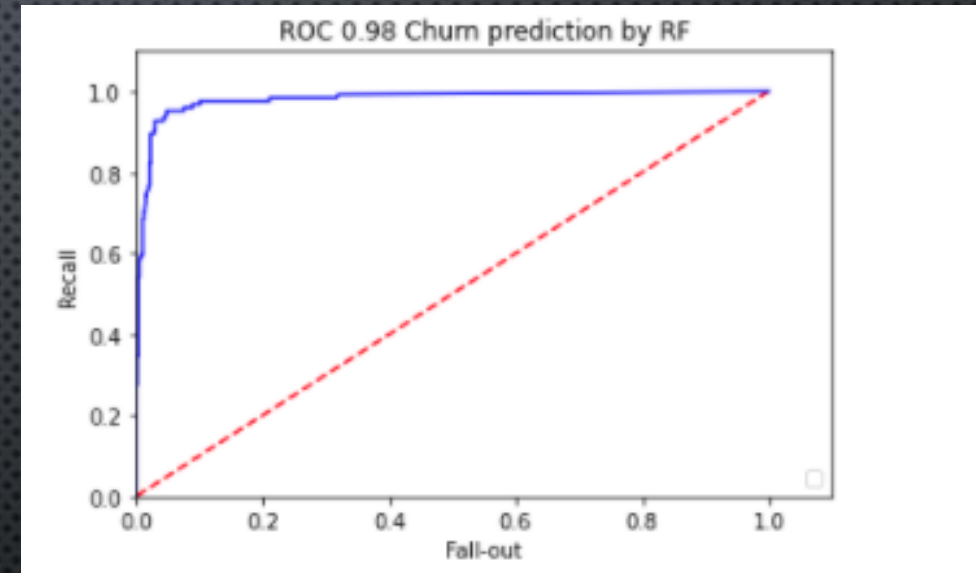
Confusion Matrix:

```
[[510  11]
 [ 22  87]]
```


RANDOM FOREST AVEC GRILLE DE RECHERCHE

En faisant évoluer les hyperparamètres à l'aide du grid search, la meilleure configuration semble être la suivante : `Min_samples_leaf = 3`, `n_estimators = 10`, `n_jobs = -1` avec une précision de 0.967

On peut apercevoir avec la courbe ROC que le modèle prédit très bien les labels et est très précis.



19:41:27

RandomForestClassifier
Fitting 5 folds for each of 1 candidates, totalling 5 fits

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

best_score_ 0.9471870834085843 with = RandomForestClassifier(min_samples_leaf=3, n_estimators=10, n_jobs=-1)
19:41:31

KNN

Avec le KNN, nous avons une précision de 0.846.

Tout comme la régression logistique, l'algorithme prédit de manière satisfaisante le label 0 qui est en majorité dans ce dataset mais on a un nombre assez conséquent de FP = 63, qui est un danger pour le besoin métier et on prédit mal les clients voulant résilier

On a une bonne précision globale et un roc score correct mais la remarque ci-dessous pourrait être un argument pour ne pas choisir ce modèle dans le cadre de ce problème

Accuracy Score:

0.846

ROC Score:

0.8350208667171459

Classification Report:

	precision	recall	f1-score	support
0	0.89	0.93	0.91	521
1	0.57	0.42	0.49	109
accuracy			0.85	630
macro avg	0.73	0.68	0.70	630
weighted avg	0.83	0.85	0.84	630

Confusion Matrix:

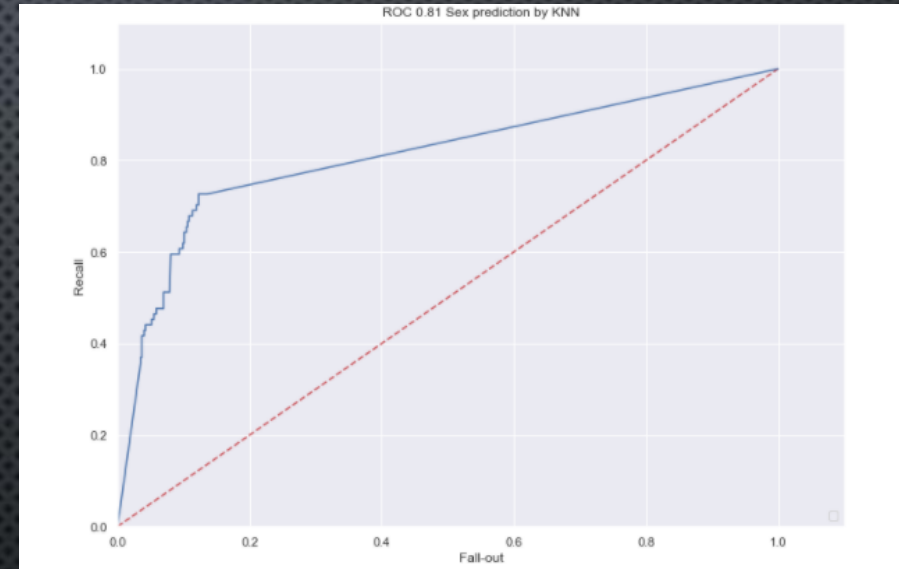
```
[[487  34]
 [ 63 46]]
```


KNN AVEC GRILLE DE RECHERCHE

En faisant évoluer les hyperparamètres à l'aide du grid search, la meilleure configuration semble être la suivante :

leaf_size = 3, n_neighbors = 2, weights = distance avec une précision de 0.967

On peut apercevoir avec la courbe ROC que le modèle prédit correctement les labels et a une bonne performance dans sa manière de classifier.



16:43:10

```
-----  
knn_grid_search  
{ 'n_neighbors': [2, 3, 5, 8, 10], 'weights': ['uniform', 'distance'], 'algorithm': ['auto', 'ball_tree', 'kd_tree',  
'brute'], 'leaf_size': [1, 3]}  
best_score_ = 0.86 with = KNeighborsClassifier(leaf_size=1, n_neighbors=2, weights='distance')  
16:43:25
```

Duration time : 15.205899953842163

GRADIENT BOOSTING

- Avec le gradient boosting, l'algorithme prédit de manière satisfaisante les labels et fait très peu d'erreur tout comme le random forest.
- C'est un algorithme que l'on pourrait facilement retenir pour le besoin de l'opérateur
- Mais elle présente un roc score plus faible que celui du random forest du à une plus faible précision de la prédiction des labels 1 qui concerne la résiliation d'un client mais cette différence est assez minime.

```
Accuracy Score:  
0.940  
ROC Score:  
0.8350208667171459
```

```
Classification Report:  
              precision    recall  f1-score   support  
  
         0         0.96      0.97      0.96         521  
         1         0.84      0.81      0.82         109  
  
    accuracy              0.94         630  
   macro avg              0.90      0.89      0.89         630  
weighted avg              0.94      0.94      0.94         630
```

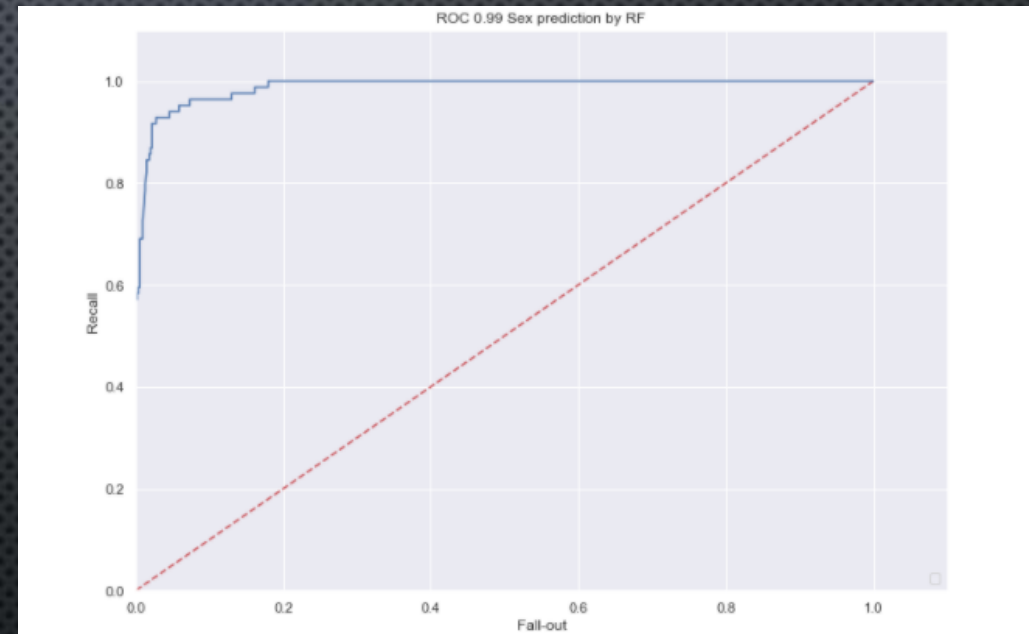
```
Confusion Matrix:  
[[504  17]  
 [ 21  88]]
```


GRADIENT BOOSTING AVEC GRILLE DE RECHERCHE

En faisant évoluer les hyperparamètres à l'aide du grid search, la meilleure configuration semble être la suivante :

'n_estimators' = 10 , 'max_depth' = None ,
'min_samples_split'=2, 'learning_rate'= 0.01
, 'loss': 'ls' avec une précision de 0.94

On peut apercevoir avec la courbe ROC que le modèle prédit très bien les labels et est très précis.



```
16:45:06
-----
clf_grid_search
{'n_estimators': 10, 'max_depth': None, 'min_samples_split': 2, 'learning_rate': 0.01, 'loss': 'ls'}
best_score_ = 0.94 with = GradientBoostingClassifier()
16:45:09

Duration time : 3.7153279781341553
```

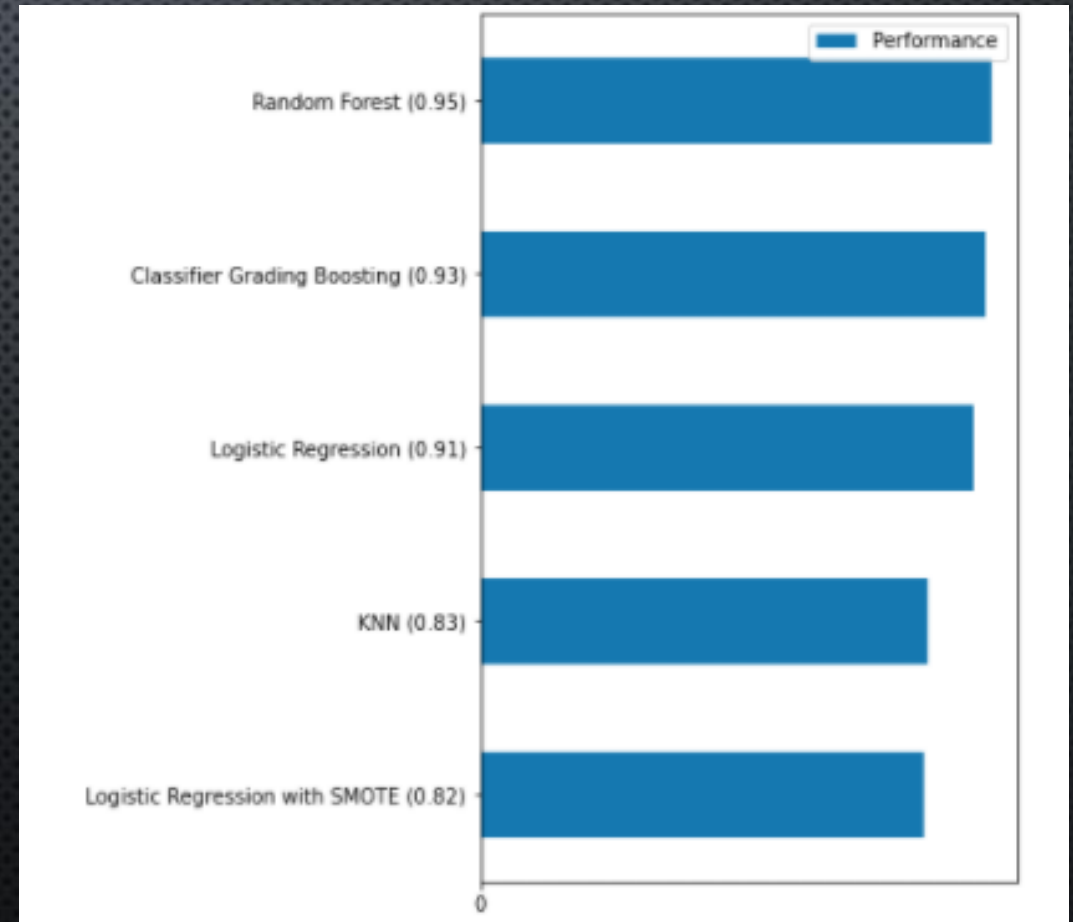
RÉSUMÉ : COMPARAISONS DES ALGORITHMES

En résumé des algorithmes utilisés dans le cadre de ce problème, tous les algorithmes ont une performance assez élevée.

Mais comme vu précédemment, certains algorithmes sont plus performants et plus prédisposés à prédire la variable cible. Le meilleur algorithme est le random forest suivi de près par le gradient boosting.

Cependant malgré la dernière place de la régression logistique avec le SMOTE, elle semble plus adaptée et prend moins de temps d'exécution que la régression logistique et le KNN pour le besoin de l'opérateur.

Le random forest serait le meilleur algorithme de par sa capacité à prédire correctement les labels pour l'opérateur et par sa rapidité d'exécution





API FLASK

L'API permet de tester des données selon les différents modèles présentés plus haut. Pour sélectionner un modèle spécifique, il suffit de préciser le type de modèle en paramètre.

API Endpoint:

- `/api/` Méthode: POST
Prédiction de la résiliation ou non d'un client en fonction d'un modèle défini dans les paramètres.

Format JSON à envoyer:

```
{  
  'values': [[8,0, 38, 0, 4370, 71, 5, 17, 3, 1, 1, 30, 69.764]],  
  'model': 'RF'  
}
```

Liste des modèles accessibles:

- LR : Logistic Regression
- LRS : Logistic Regression with SMOTE
- RF : Random Forest
- KNN : K-Nearest Neighbours
- CLF : Gradient Boosting