



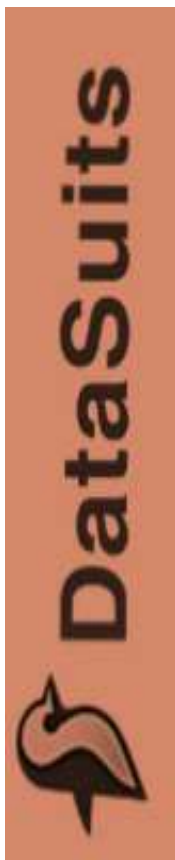
# Formation Data Analyst

## Projet Final :

**Conception de modèles prédictifs pour  
l'approbation d'un crédit bancaire.**

Du 18 aout au 24 octobre 2025

Roger Mamaty



# Table des Matières

1- <a href="#">Présentation</a> .....	3
1-1 <a href="#">Jeu de données</a> .....	3
1-2 <a href="#">Catégorie des données</a> .....	3
2- Analyse exploratoire des données sous Python.....	5
2-1 <a href="#">Intégration et conversion des données</a> .....	5
2-2 <a href="#">Traitement des valeurs aberrantes</a> .....	5
2-3 <a href="#">Matrice de Corrélation</a> .....	6
3- <a href="#">Régression Logistique</a> .....	7
3-1 <a href="#">Évaluation du modèle</a> .....	7
3-2 <a href="#">matrice de confusion</a> .....	7
4- Forêts Aléatoire.....	8
4-1 <a href="#">Évaluation du modèle</a> .....	8
4-2 <a href="#">matrice de confusion</a> .....	8
4-3 <a href="#">visualisation de l'arbre a 2 niveaux</a> .....	9
5- <a href="#">Comparaison Régression Logistique / Forêts Aléatoire</a> .....	10
6- <a href="#">Conclusion</a> .....	11
7- <a href="#">Annexe : Présentation des données sous PowerPoint</a> .....	12
7-1 <a href="#">Analyse Démographique</a> .....	12
7-2 <a href="#">Analyse Emplois et Revenus</a> .....	13
7-3 <a href="#">Analyse Financière</a> .....	14
7-4 <a href="#">Analyse Paramètres Crédit</a> .....	15

# 1- Présentation

Dans le cadre du rapport final de la formation de « Data Analyst » auprès de **DataSuits**, j'ai choisi de travailler sur des modèles prédictifs pour obtention du crédit : oui / non. Le jeu de données retenue est Loan.csv (<https://www.kaggle.com/datasets/lorenzozoppelletto/financial-risk-for-loan-approval/data>).

L'objectif est de développer des modèles prédictifs afin de déterminer le résultat binaire de l'approbation d'un prêt, indiquant si un demandeur est susceptible d'être approuvé ou refusé pour un prêt.

## 1-1 Jeu de données

Le fichier **Loan.csv** contient 36 colonnes et 20 000 enregistrements. Une phase de traduction est réalisée sur les colonnes origines afin de simplifier l'usage de ces colonnes dans les différents processus d'analyse.

Les données comprennent de nombreuses informations, nous allons les présenter en les regroupant par type d'information : *Temporelles*, *Démographiques*, *Emplois et Revenus*, *logements*, *Financières*, *Charges et endettement*, *Crédit*, *Prêt*, *Risque*.

## 1-2 Catégorie des données

### Variables temporelles

Variable	Description
Date_demande	Date de soumission de la demande de prêt.

### Variables liées au logement

Variable	Description
Statut_proprietaire	Statut résidentiel (locataire, propriétaire...).

### Variables démographiques

Variable	Description
Age	Âge du demandeur.
Etat_civil	État matrimonial (célibataire, marié, divorcé...).
Nbre_Personnes_a_Charges	Nombre de personnes à charge.
Niveau_Education	Plus haut niveau d'éducation atteint.

### Variables liées à l'emploi et aux revenus

Variable	Description
Statut_Emploi	Situation professionnelle (CDI, CDD, indépendant...).
Experience	Années d'expérience professionnelle.
Duree_Emploi_Actuel	Durée dans l'emploi actuel.
Revenu_annuel	Revenu annuel.
Revenu_Mensuel	Revenu mensuel (souvent dérivé de AnnualIncome).

### Variables financières générales

Variable	Description
<b>Solde_Compte_Courant</b>	Solde du compte courant.
<b>Solde_Compte_Epargne</b>	Solde du compte épargne.
<b>Total_Actifs_Detenus</b>	Valeur totale des actifs détenus.
<b>Total_Dettes-Dus</b>	Montant total des dettes.
<b>Valeur_nette</b>	Valeur nette (TotalAssets - TotalLiabilities).

### Variables de charges et endettement

Variable	Description
<b>Palements_mensuels_dettes</b>	Palements mensuels de dettes existantes.
<b>Ratio_Dette_Revenu</b>	Ratio dettes / revenu.
<b>Ratio_DettePret_Revenu</b>	Ratio total dettes / revenu en tenant compte du prêt demandé.
<b>Historique_paiement_charges</b>	Historique de paiement des charges/abonnements (eau, électricité...).

### Variables liées au crédit

Variable	Description
<b>CreditScore</b>	Score de crédit mesurant la solvabilité.
<b>Duree_Historique_Credit</b>	Durée totale de l'historique de crédit.
<b>Historique_comportement_paiement</b>	Historique des paiements passés (ponctualité, retards...).
<b>Taux_Utilisation_carte_de_credit</b>	Pourcentage d'utilisation des cartes de crédit.
<b>Nombre_lignes_credit_actifs</b>	Nombre de lignes de crédit actives.
<b>Nombre_verifications_credits</b>	Nombre de vérifications de crédit récent.
<b>Historique_Faillite</b>	Historique de faillites.
<b>Precedent_Defaults_de_paiement</b>	Historique de défauts de paiement sur prêts précédents.

### Variables liées au prêt

Variable	Description
<b>Montant_Pret</b>	Montant du prêt demandé.
<b>Duree_Pret</b>	Durée de remboursement du prêt.
<b>Motif_Pret</b>	Objet du prêt (voiture, travaux...).
<b>Taux_Interet_de_base</b>	Taux d'intérêt de base du marché.
<b>Taux_interet_Obtenu</b>	Taux d'intérêt appliqué au demandeur.
<b>Mensualite_a_payer_pret</b>	Mensualité à payer pour le prêt.
<b>Pret_Approuve</b>	Statut d'approbation du prêt (oui/non).

### Variables de risque

Variable	Description
<b>RiskScore</b>	Score global de risque calculé par l'institution.

## 2- Analyse exploratoire des données sous Python

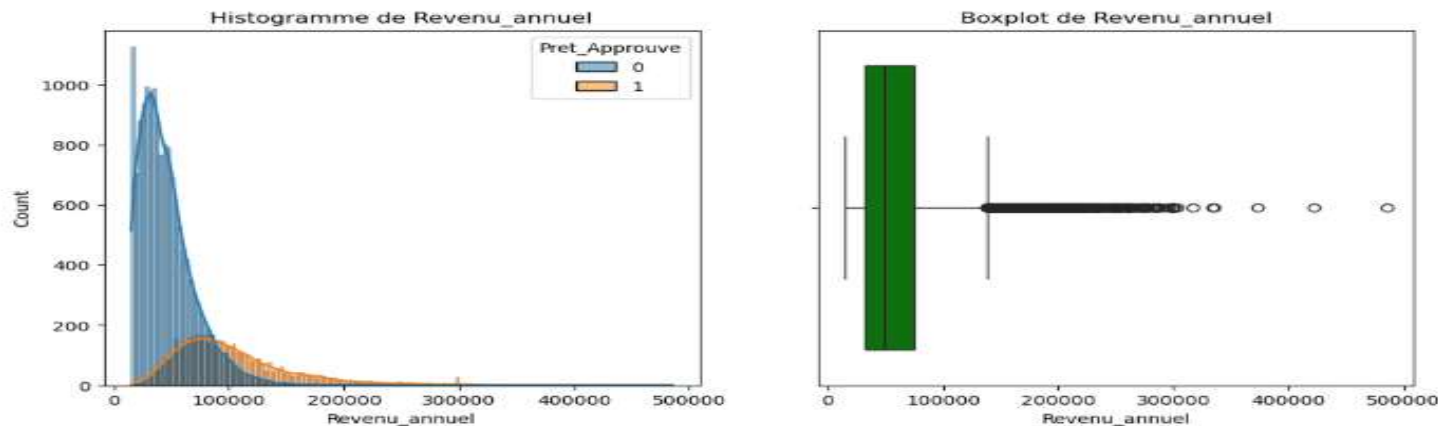
### 2-1 Intégration et conversion des données

Le fichier Loan.CSV est intégré dans Python. Après intégration des données, nous procédons au formatage des données :

- La colonne date est transformé en type Date
- 3 colonnes sont transformées en type Catégorie.
- Pas de doublons détectés dans les données du datagramme.
- De nouvelles création de colonnes seront créées afin de regrouper les colonnes à valeurs numériques en libellé du type *faible*, *moyen*, *élevé*. Qui vont regrouper des intervalles de valeur. Ces colonnes vont nous permettre un regroupement sur des visuels de type PowerPoint ou Excel.

### 2-2 Traitement des valeurs aberrantes

Nous allons utiliser l'affichage des Boxplot des colonnes pour visualiser d'éventuelles valeurs aberrantes. Dans l'exemple ci-dessous, on constate de possibles valeurs Outliers et ou aberrantes. Après réduction de certaines valeurs Outliers, nous avons un jeu de données de 19830 lignes et 36 colonnes.

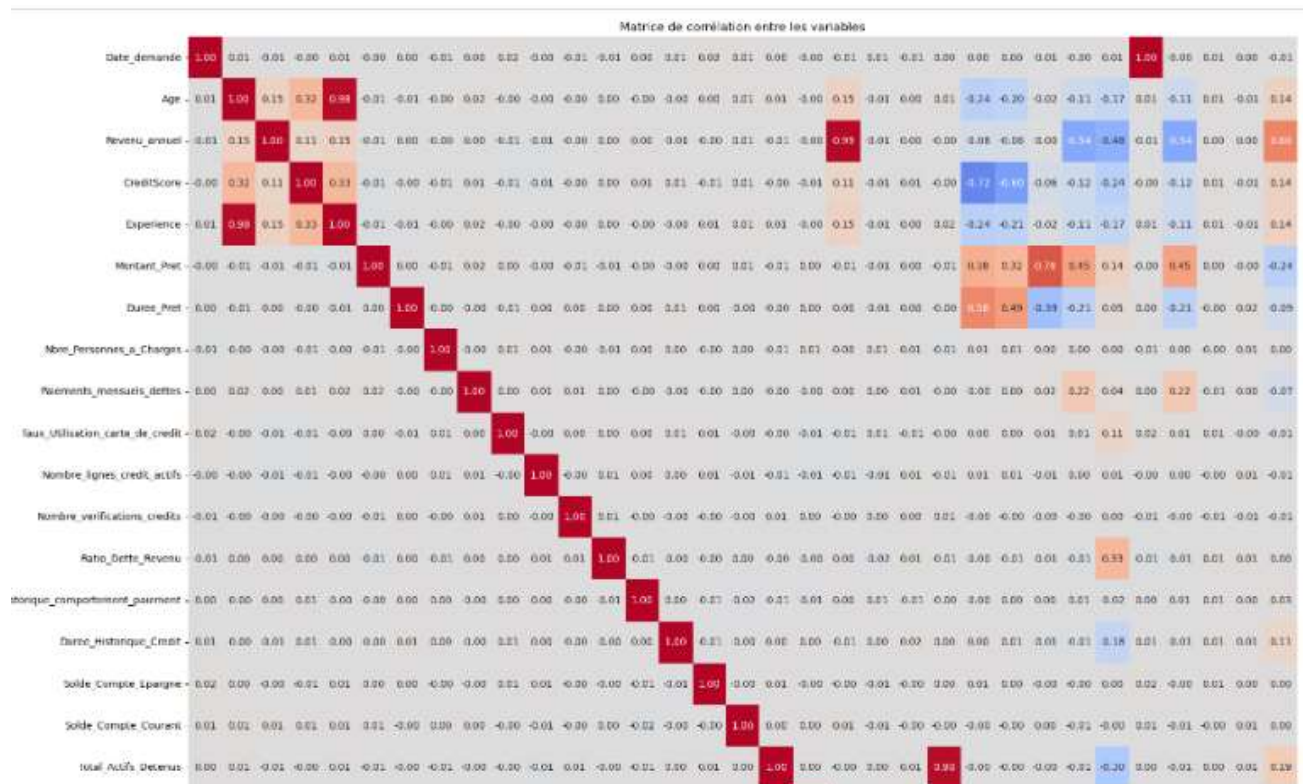


## 2-3 Matrice de Corrélation

La matrice de corrélation va nous permettre de voir :

- Les features qui ont le plus d'impact (positif ou négatif) sur l'obtention du crédit (Pret\_Approuve).
- Exclure les features ayant une forte multi colinéarité.
- Exclure les features ayant un impact nul sur l'obtention du crédit (Pret\_Approuve).

Extrait de la matrice  
de corrélation



## 3- Régression Logistique

Notre premier modèle à évaluer est la « Régression Logistique ». Ce modèle va nous permettre de détecter les « mosts features » lié à ce modèle, et de réaliser un modèle de prédiction l'approbation ou non du crédit bancaire.

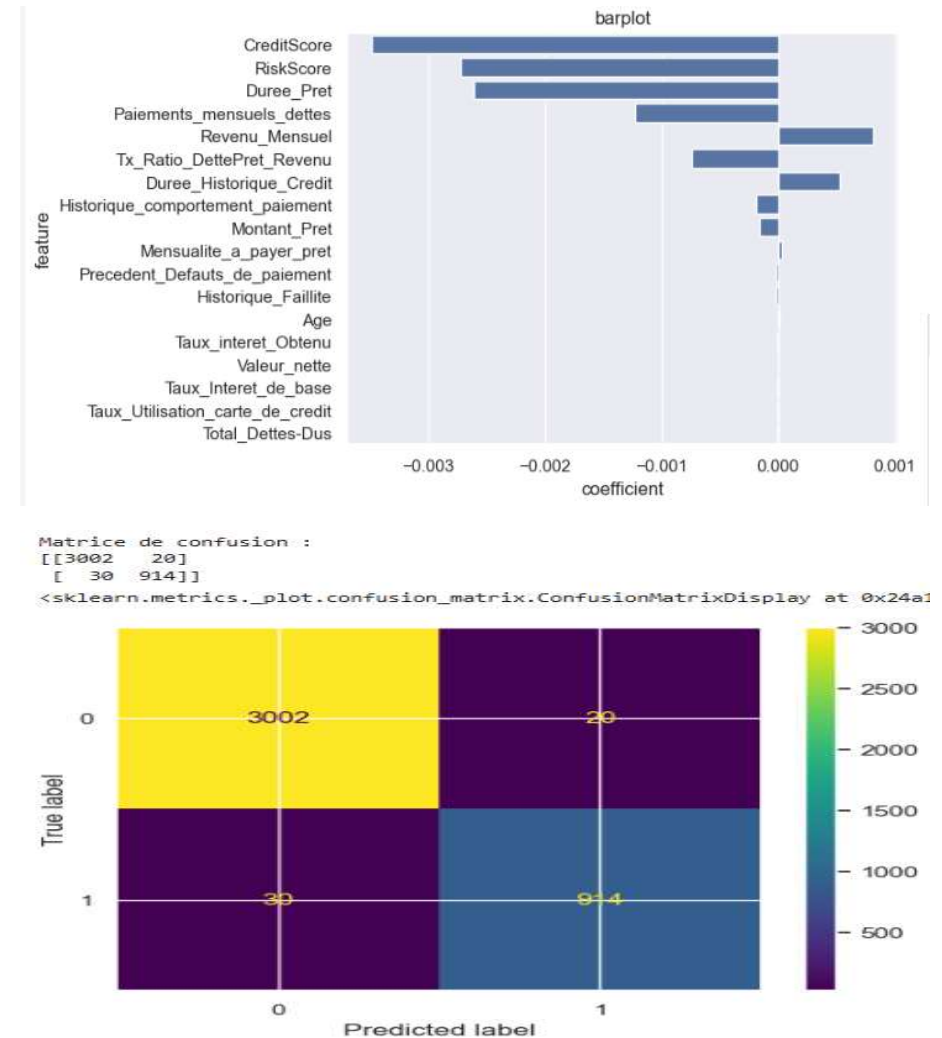
### 3-1 Évaluation du modèle

Le résultat du modèle de régression logistique nous donne les informations suivantes :

- **Précision du modèle en %** : **98.74%**
- Nombre de prédictions positifs bien effectuées : 97.86%
- Nombre de positifs (approuvé) bien prédit : 96.82%

### 3-2 matrice de confusion

- Vrais négatifs (TN) = 3002.
- Faux positifs (FP) = 20.
- Faux négatifs (FN) = 30.
- Vrais positifs (TP) = 914.



## 4- Forets Aléatoire

Nous allons utiliser un autre modèle, « la Forêt Aléatoire » et évaluer sa performance à générer un modèle de prédiction d'acceptation du crédit

### 4-1 Évaluation du modèle

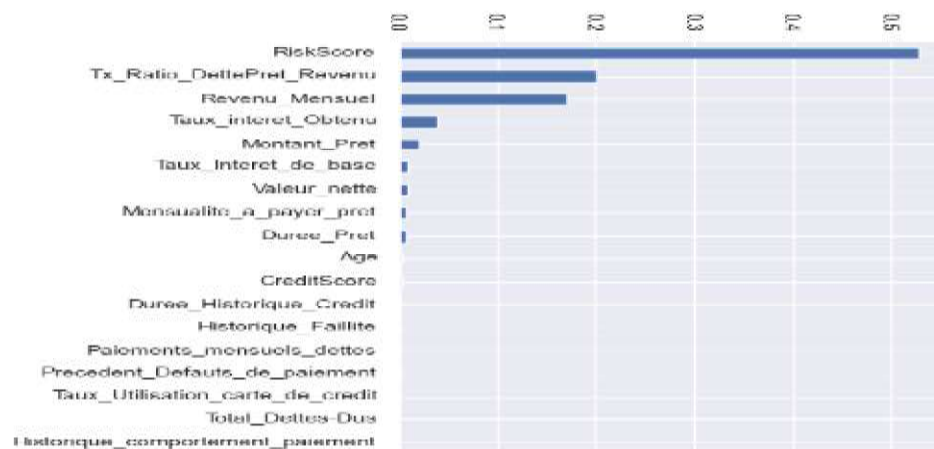
Le résultat du modèle de régression logistique nous donne les informations suivantes :

- **Précision du modèle en %** : **98.39%**
- Nombre de prédictions positifs bien effectuées : 97.31%
- Nombre de positifs (approuvé) bien prédit : 95.87%

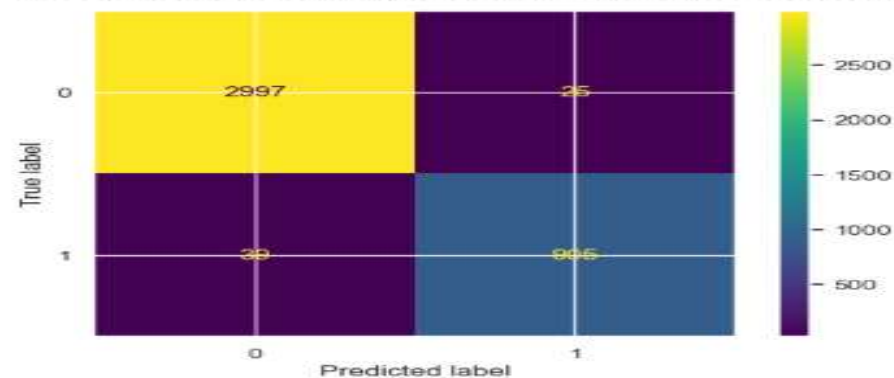
### 4-2 matrice de confusion

- Vrais négatifs (TN) = 2997.
- Faux positifs (FP) = 25.
- Faux négatifs (FN) = 39.
- Vrais positifs (TP) = 905.

INTERPRÉTATION du Most Important features: FORÊT ALÉATOIRE



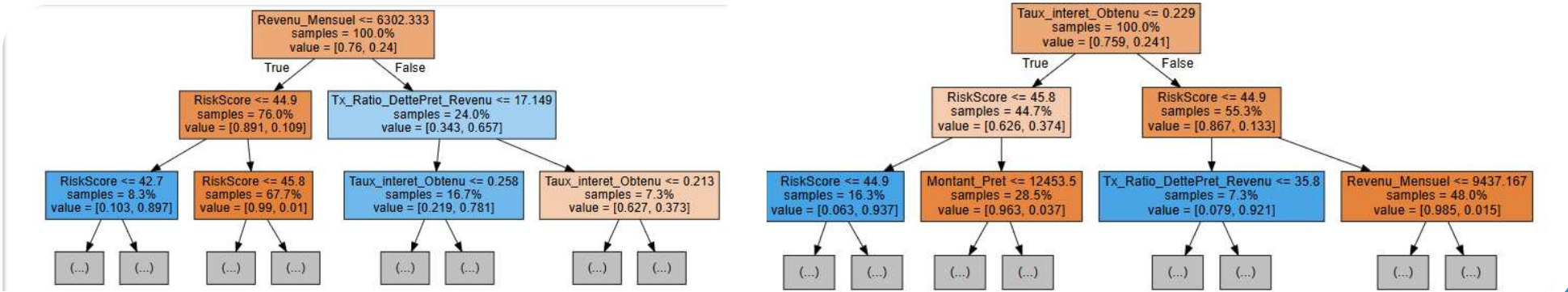
Matrice de confusion Forêt Aléatoire:  
[[ 2997 25]  
[ 39 905]]  
<sklearn.metrics.\_plot\_confusion\_matrix.ConfusionMatrixDisplay at 0x1bd85292358>





# 4- Forets Aléatoire

## 4-3 visualisation de l'arbre a 2 niveaux



## 5- Comparaison Régression Logistique / Forets Aléatoire

### Most feature Régression Logistique

1-CreditScore
2-RiskScore
3-Duree_Pret
4-Paiements_mensuels_dettes
5-Revenu_Mensuel
6-Tx_Ratio_DettePret_Revenu
7- Duree_Historique_Credit
8-Historique_comportement_paiement
9-Montant_Pret
10- Mensualite_a_payer_pret

### Comparaison Régression Logistique / Foret Aléatoire

Libellé	Régression Logistique	Forets Aleatoire
<b>Accuracy</b>	<b>98,74%</b>	<b>98,39%</b>
<b>Precision</b>	<b>97,86%</b>	<b>97,31%</b>
<b>Recall</b>	<b>96,82%</b>	<b>95,87%</b>

### Most Feature Forets Aléatoire

1-RiskScore
2-Tx_Ratio_DettePret_Revenu
3-Revenu_Mensuel
4-Taux_interet_Obtenu
5-Montant_Pret
6-Taux_interet_de_base
7-Valeur_nette
8- Mensualite_a_payer_pret
9-Duree_Pret

### Interprétation :

- **accuracy** calcule la précision du modèle, c'est-à-dire le pourcentage de prédictions correctes par rapport à l'ensemble de test.
- **precision** est utilisée pour calculer la précision du modèle concernant la proportion de vrais positifs (client ayant le pret approuvé) parmi les predictions des clients positifs. c'est le nombre de positifs bien prédit (Vrai Positif) divisé par l'ensemble des positifs prédit (Vrai Positif + Faux Positif).
- **recall** est utilisée pour calculer le rappel du classificateur, c'est le nombre de positifs bien prédit (Vrai Positif) divisé par l'ensemble des positifs (Vrai Positif + Faux Négatif).

## 6- Conclusion

À la suite de l'analyse Python, les deux modèles prédictifs présentent des critères impactant différents.

Concernant la « **Régression Logistique** », ses 3 premiers critères impactant sont **CreditScore**, **RiskScore** et la **durée\_prêt**. Tandis que la « **Forêt Aléatoire** », ses 3 premiers critères impactant sont **RiskScore**, **Tx\_Ratio\_DettePret\_Revenu** et **Revenu\_Mensuel**.

La « **Régression Logistique** » privilégie les features qui sont des indicateurs calculés par la banque, alors que la « **Forêt Aléatoire** » privilégie les informations du demandeur (revenu et dettes).

La précision générale (*accuracy*), est élevé pour les deux modèles, plus de 98%, « la **Régression Logistique** » a le taux le plus élevé de réussite globale (98,74% contre 98,39%).

Concernant l'efficacité à trouver les vrais positifs (*precision*), %, « la **Régression Logistique** » (97,86%) a un taux supérieur à celui de la « **Forêts Aléatoire** » (97,31%).

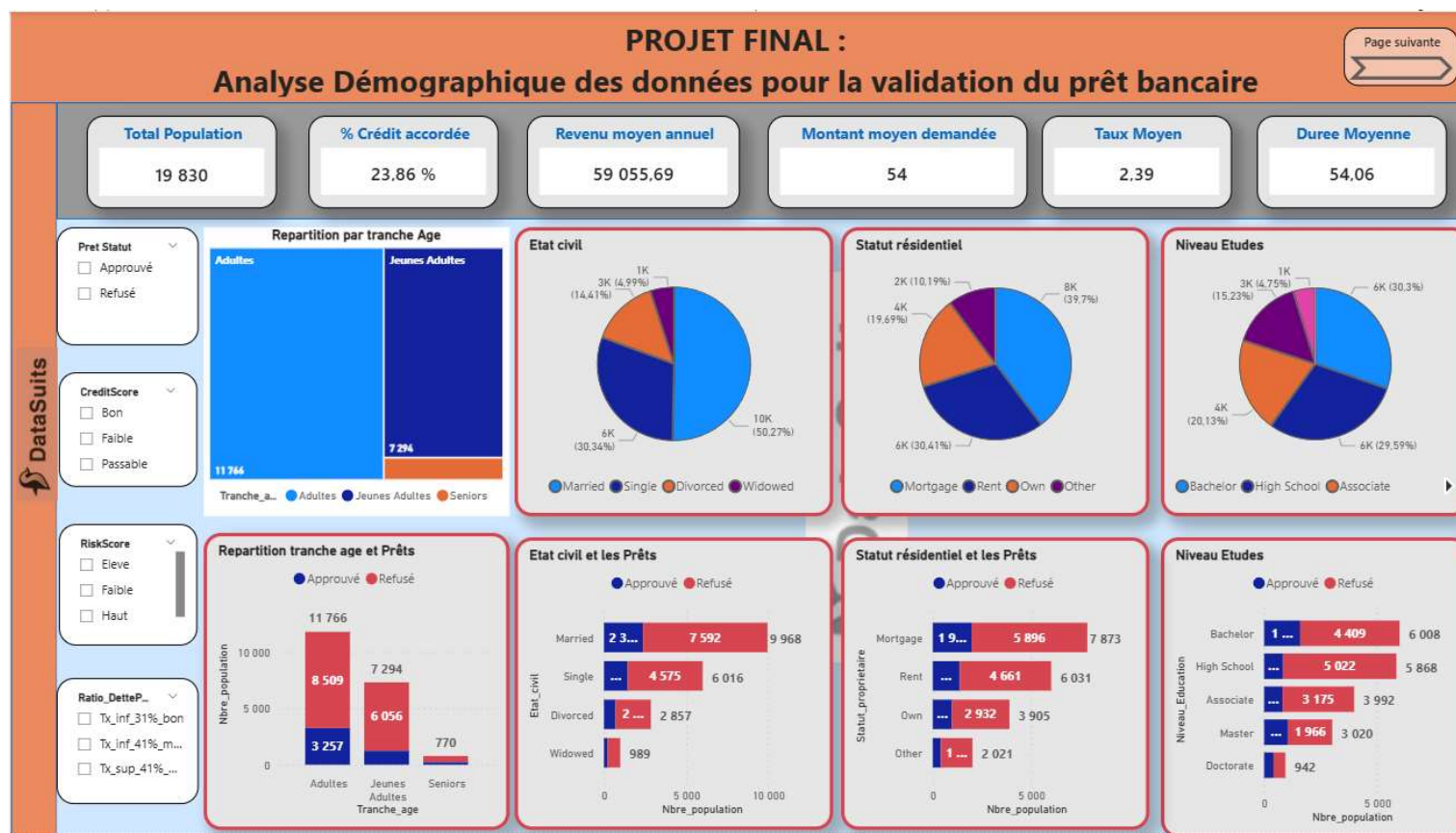
Finalement, la « **Régression Logistique** » est plus adaptée pour prédire au mieux un profil client peut voir sa demande de crédit acceptée ou non. Elle se passe en priorité sur les indicateurs établis par la banque à partir des données du demandeur.

La « **Forêt Aléatoire** » utilise au mieux les informations usuels du client pour la prédiction, malgré son taux de réussite à plus de 98%, elle est moins efficace.

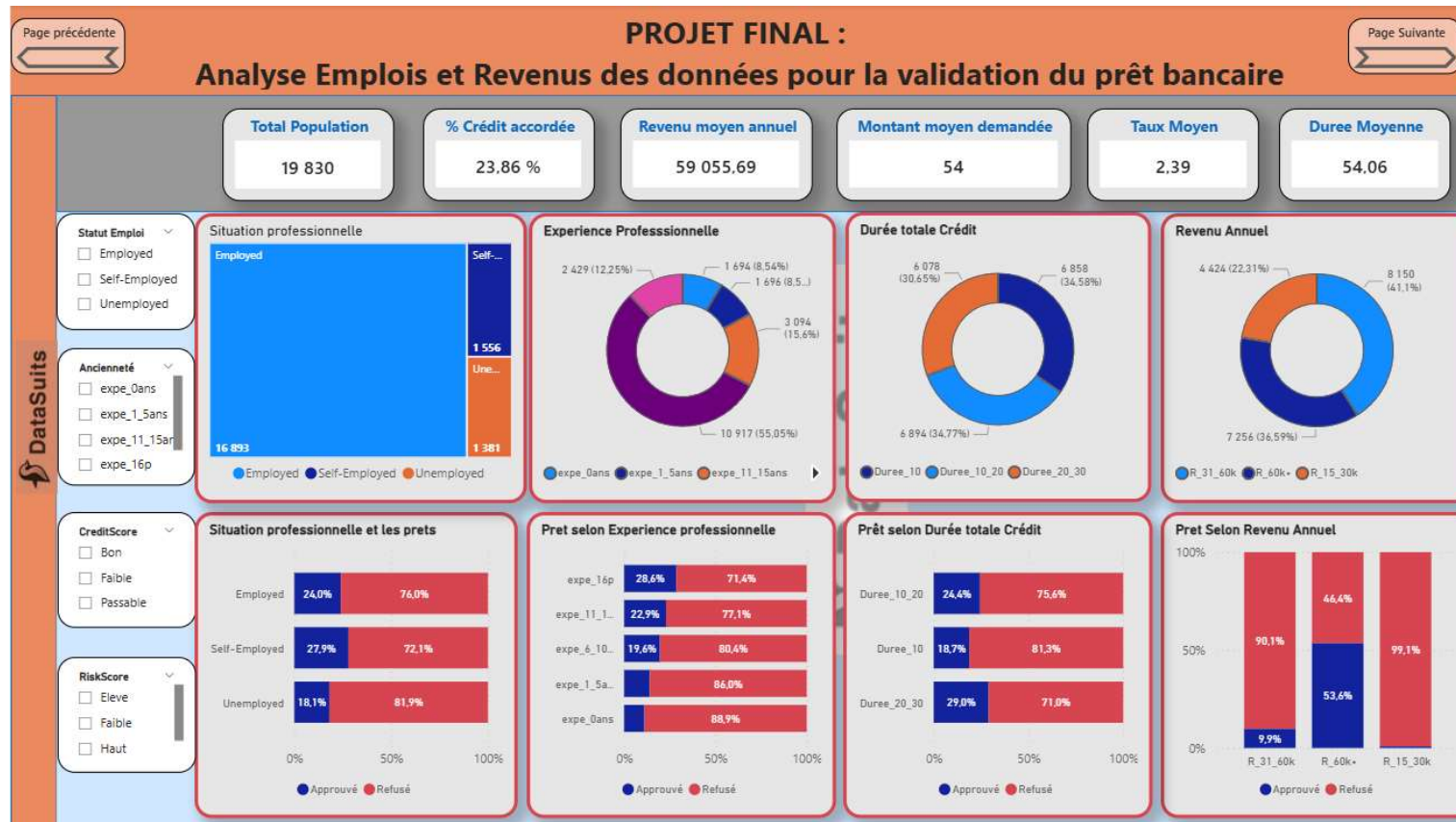
**Nb:** Les documents Notebook Python et PowerBy sont consultables via ce lien : [GitHub - Rogermamaty/projet-final: Projet final - Formation Data Analyst](#)

## 7- Annexe : Présentation des données sous PowerPoint

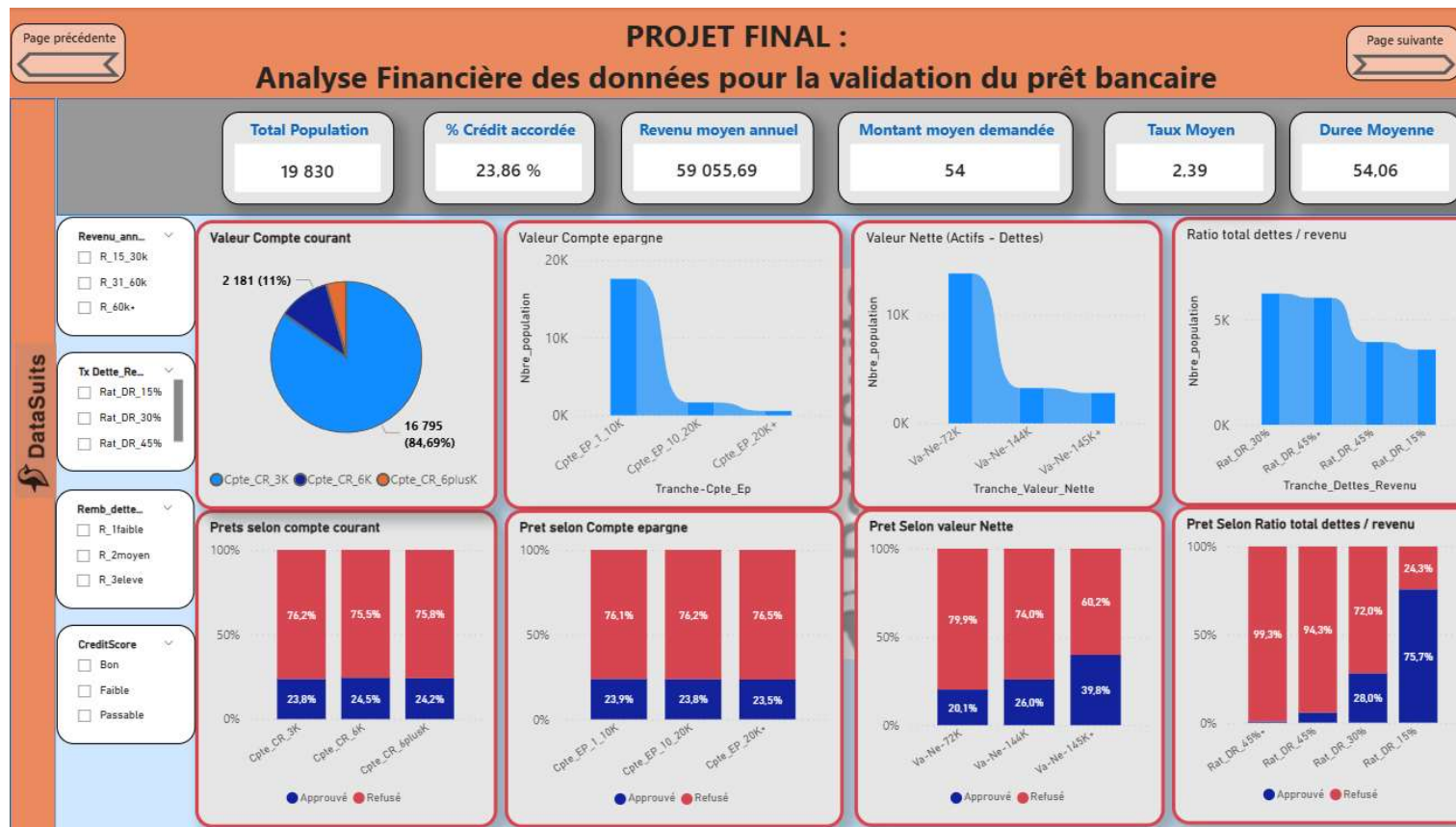
### 7-1 Analyse Démographique



## 7-2 Analyse Emplois et Revenus



## 7-3 Analyse Financière





## 7-4 Analyse Financière

