



Projet Scoring - Filière AMOA INE2

Présenté à :

INSTITUT NATIONAL DES POSTES ET TELECOMMUNICATIONS

Sujet

Prévision des difficultés financières des clients d'une banque

Réalisé par :

Achraf KHAMRI
Khalil LAMHARCHI
Ziyad GHANEME
Abdessamad ABID
Walid ELASS

Encadré par :

Pr. Zineb EL AKKAOUI

Année universitaire :

2020/2021

Résumé

Dans le cadre de nos études au sein de l'INPT en deuxième année, nous - étudiants de la filière AMOA - étions amenés à réaliser un projet où nous devrions employer l'ensemble de nos acquis.

Ce projet s'inscrit dans le thème de l'analyse de données avancée et le Scoring. Divisée en plusieurs groupes, la filière AMOA était amenée à travailler afin de proposer une solution adéquate au sujet. Ce document contient les résultats du travail de notre groupe.

Durant la réalisation de ce projet, nous avons parcouru toutes les phases de le l'anaylse de données a savoir le prétraitement, la création de modèle de prévision et l'évaluation de la performance en s'appuyant sur diverses technologies, principalement Python.

Liste des abréviations

- **AMOA** - Assistance à maîtrise d'ouvrage
- **ACP** - Analyse des Composantes Principales
- **LDA** - Linear Discriminant Analysis
- **QDA** - Quadratic Discriminant analysis
- **SQL** - Structured Query Language
- **TPR** - True Positive Rate
- **FPR** - False Positive Rate
- **TNR** - True Negative Rate
- **FNR** - False Negative Rate
- **ROC** - Receiver Operating Characteristic

Liste des figures

1	Aperçu des données	6
2	Logo de Pandas	8
3	Logo de Microsoft SQL Server	8
4	Requête taux de défaut	9
5	Boîtes à moustaches	9
6	Tableau des intervalles de données valides	10
7	Matrice de nullité	10
8	Valeurs complètes par colonne	11
9	Requête SQL de l'équilibrage	11
10	Comparaison entre les 2 classes de clients	12
11	Séparation des groupes en LDA	13
12	Régression Logistique	14
13	Scores de précision des modèles	14
14	Matrice de confusion LDA	15
15	Matrice de confusion QDA	15
16	Matrice de confusion Logistic Regression	15
17	Tableau des pouvoirs de prédiction	15
18	Courbe ROC de Logistic Regression	16
19	Matrice de corrélation	17

Table des matières

Résumé	1
Liste des abréviations	2
Liste des figures	3
1 Contexte du projet	5
1.1 Introduction	5
1.2 Cadre du projet	5
1.2.1 Présentation de l'existant	5
1.2.2 Problématique	6
1.2.3 Solution	6
2 Techniques et environnement logiciel	7
2.1 Introduction	7
2.2 Techniques	7
2.2.1 Nettoyage de données	7
2.2.2 Régression linéaire Stochastique	7
2.2.3 Analyse discriminante	7
2.3 Environnement logiciel	8
2.3.1 Pandas	8
2.3.2 Microsoft SQL Server	8
3 Mise en oeuvre de la solution	9
3.1 Phase de prétraitement	9
3.1.1 Préparation des données	9
3.2 Modèle de prévision	12
3.2.1 Analyse factorielle discriminante	12
3.2.2 Linear & Quadratic Discriminant Analysis	13
3.2.3 Logistic Regression	13
3.3 Phase d'évaluation et règle de décision retenue	14
3.3.1 Goodness of fit	14
3.3.2 Pouvoir de prédiction	15
3.3.3 Conclusion	16
4 Références bibliographiques	18
5 Annexe	19

1 Contexte du projet

1.1 Introduction

Le scoring est une technique qui permet d'affecter un score à un client ou prospect. Le score obtenu traduit généralement la probabilité qu'un individu réponde à une sollicitation marketing ou appartienne à la cible recherchée. Il mesure donc l'appétence pour l'offre potentielle. Le score peut également exprimer la valeur potentielle d'un client ou prospect ou l'espérance mathématique de gain qui peut lui être associée, voire un risque dans le cadre du scoring bancaire.

1.2 Cadre du projet

1.2.1 Présentation de l'existant

Nous avons à notre disposition une base de données regroupant les informations suivantes sur 150.000 emprunteurs :

- **SeriousDlqin2yrs**
La personne était en désordre de paiement de 90 jours ou pire.
- **RevolvingUtilizationOfUnsecuredLines**
Solde total sur les cartes de crédit et les marges de crédit personnelles à l'exception de l'immobilier et sans dette à tempérament comme les prêts automobiles divisé par la somme des limites de crédit.
- **Age**
L'âge de l'emprunteur en années.
- **NumberOfTime30-59DaysPastDueNotWorse**
Nombre de fois où l'emprunteur est en souffrance depuis 30 à 59 jours, mais pas pire au cours des 2 dernières années.
- **DebtRatio**
Paiements mensuels de la dette, pension alimentaire, frais de subsistance divisés par le revenu brut mensuel.
- **MonthlyIncome**
Revenu mensuel.
- **NumberOfOpenCreditLinesAndLoans**
Nombre de prêts ouverts (à tempérament comme prêt automobile ou hypothécaire) et marges de crédit (par exemple, cartes de crédit).
- **NumberOfTimes90DaysLate**
Nombre de fois où l'emprunteur est en souffrance depuis 90 jours au cours des 2 dernières années.
- **NumberRealEstateLoansOrLines**
Nombre de prêts hypothécaires et immobiliers, y compris les marges de crédit hypothécaire.
- **NumberOfTime60-89DaysPastDueNotWorse**
Nombre de fois où l'emprunteur est en souffrance depuis 60 à 89 jours, mais pas pire au cours des 2 dernières années.
- **NumberOfDependents**
Nombre de personnes à charge dans la famille excluant l'emprunteur (conjoint, enfants, etc.)

Voici un aperçu des données en Excel

FIGURE 1 – Aperçu des données

1.2.2 Problématique

Identifier les clients (emprunteurs) qui auraient des difficultés financières dans les deux années à venir afin de leur proposer des offres adaptées.

1.2.3 Solution

Construire des **modèles de prévision de score**. Ce score servira à attribuer **une classe** à chaque client. Les modèles construits vont exploiter les données des fichier *ScoringTraining.csv* et *ScoringTest.csv* pour l'apprentissage et la validation, ainsi que pour l'évaluation des modèles afin d'en choisir le plus performant.

2 Techniques et environnement logiciel

2.1 Introduction

Ce chapitre sera dédié à la présentation des techniques et logiciels sur lesquelles on s'est basé pour réussir l'implémentation de notre solution.

2.2 Techniques

2.2.1 Nettoyage de données

Les données dans leur forme brute peuvent avoir plusieurs erreurs et valeurs manquantes. Pour amener le modèle de prévision à prendre des décisions correctes, il est nécessaire que les informations soient fiables et complètes. Nous allons nous servir de représentations graphiques de données statistiques, notamment les boîtes à moustaches dans le but de repérer les valeurs aberrantes de chaque variable. En ce qui concerne les valeurs manquantes, nous allons les remplir en utilisant une régression linéaire multiple stochastique.

2.2.2 Régression linéaire Stochastique

C'est la solution que nous proposons pour traiter les valeurs manquantes, elle consiste à remplacer ces données par les valeurs prédites par un modèle de régression et répéter ce processus pour chaque variable. Un inconvénient majeur de cette méthode est que nous réduisons la variabilité de la variable imputée. En d'autres termes, puisque nous remplaçons les données manquantes par des sorties de régression, les valeurs prédites se trouvent au long de l'hyperplan de régression où les variables auraient en fait contenu du bruit / biais. Pour ajouter de l'incertitude aux valeurs des variables imputées, nous pouvons ajouter du bruit normalement distribué avec une moyenne de zéro et la variance égale à l'erreur-type des estimations de régression.[1]

2.2.3 Analyse discriminante

C'est un ensemble de méthodes statistiques de réduction de dimension qui vise à prédire l'appartenance d'un individu à des groupes prédéfinis à partir d'une série de variables prédictives. Le principe consiste à projeter les individus dans une direction permettant de mettre en évidence les groupes. À cette fin, Il faut privilégier la variance interclasse au détriment de la variance intraclasse considérée comme due au bruit. Dans la mise en œuvre de notre solution nous allons implémenter une analyse factorielle discriminante (AFD), une Analyse discriminante linéaire (LDA) et une analyse discriminante quadratique (QDA).

2.3 Environnement logiciel

2.3.1 Pandas



FIGURE 2 – Logo de Pandas

Pandas est une bibliothèque spécialisée dans la manipulation de données. Cette bibliothèque contient un ensemble de fonctions optimisées pour la gestion des Big Data. Elle permet de créer et d'exporter des tableaux de données à partir de fichiers texte (séparateurs, .csv, format fixe, compressé), binaires (HDF5 avec Pytable), HTML, XML, JSON, MongoDB, SQL... [2]

2.3.2 Microsoft SQL Server



FIGURE 3 – Logo de Microsoft SQL Server

SQL(sigle de Structured Query Language, en français langage de requête structurée) est un langage informatique normalisé servant à exploiter des bases de données relationnelles. La partie langage de manipulation des données de SQL permet de rechercher, d'ajouter, de modifier ou de supprimer des données dans les bases de données relationnelles.[3]

3 Mise en oeuvre de la solution

3.1 Phase de prétraitement

3.1.1 Préparation des données

Proportion de défauts

Le taux de défaut est un ratio calculant le pourcentage de non-paiement des échéances d'un emprunt[4] et se calcule à partir de la formule suivante [5] :

$$p = \text{NombreDePrtsEnRetardOuEnDfautDePaiementDepuisPlusDeuxMois} \div \text{NombreDePrtsEnCours}$$

Ces données sont disponibles dans les colonnes **NumberOfOpenCreditLinesAndLoans**, **NumberOfTimes90DaysLate**, **NumberRealEstateLoansOrLines**, **NumberOfTime60-89DaysPastDueNotWorse**. La requête SQL suivante permet d'effectuer ce calcul

```
SELECT ( Sum(numberoftimes90dayslate)
        + Sum(numberoftime60-89dayspastduenotworse) ) / ( Sum(numberofopencreditlinesandloans)
        + Sum(numberrealestateloansorlines) )
FROM   dbo.scoringtraining
```

FIGURE 4 – Requête taux de défaut

Résultat : 5.35%

Données extrêmes

Une donnée extrême ou aberrante est une valeur ou une observation qui est distante des autres observations effectuées sur le même phénomène. Afin de détecter ces valeurs, on se sert des **boîtes à moustache** générées pour chaque variable. Le code permettant d'obtenir ces figures est disponible en annexe.

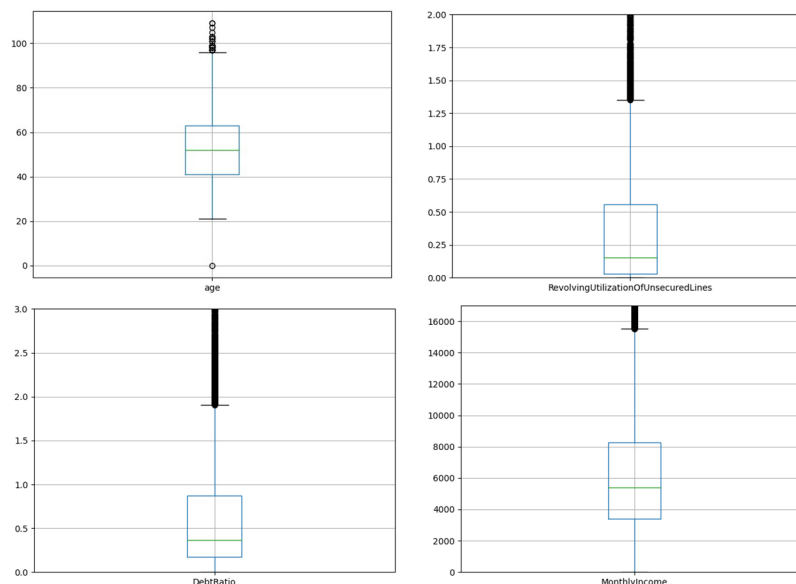


FIGURE 5 – Boîtes à moustaches

En se basant sur les graphes. Nous avons décidé de conserver les intervalles de données valides suivants :

Variable	Intervalle
RevolvingUtilizationOfUnsecuredLines	0-1.4
Age	20-98
NumberOfTime3059DaysPastDueNotWorse*	NombreCredits*24
DebtRatio	0-1.4
RevolvingUtilizationOfUnsecuredLines	0-1.9
MonthlyIncome	100-15800
NumberOfOpenCreditLinesAndLoans	0-20
NumberOfTimes90DaysLate*	NombreCredits*24
NumberRealEstateLoansOrLines	0-5
NumberOfTime6089DaysPastDueNotWorse*	NombreCredits*24
NumberOfDependents	0-4

***Ces variables ne doivent pas dépasser la valeur : (nombre de crédits) *24 puisque les données s'étalent en durée sur 2 années.**

FIGURE 6 – Tableau des intervalles de données valides

Valeurs manquantes

Le package **missingno** de python fournit un ensemble d'outils de visualisations de données manquantes flexibles et faciles à utiliser qui permettent d'obtenir un résumé visuel rapide de l'exhaustivité (ou de l'absence de celle-ci) d'un ensemble de données. [6]

L'application de ces outils sur notre base de données donne les résultats suivants :

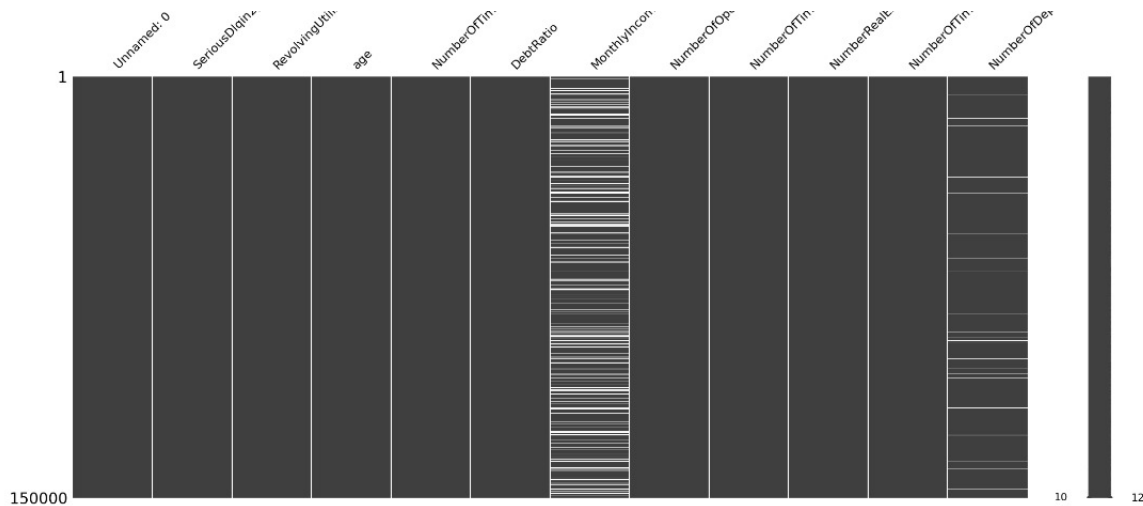


FIGURE 7 – Matrice de nullité

La matrice de nullité permet de visualiser la distribution des données manquantes. Remarquons donc que les variables concernées par ce problème sont **MonthlyIncome** et **NumberOfDependents** dans une moindre mesure. Nos remarques sont confirmées par la visualisation de la nullité par colonne

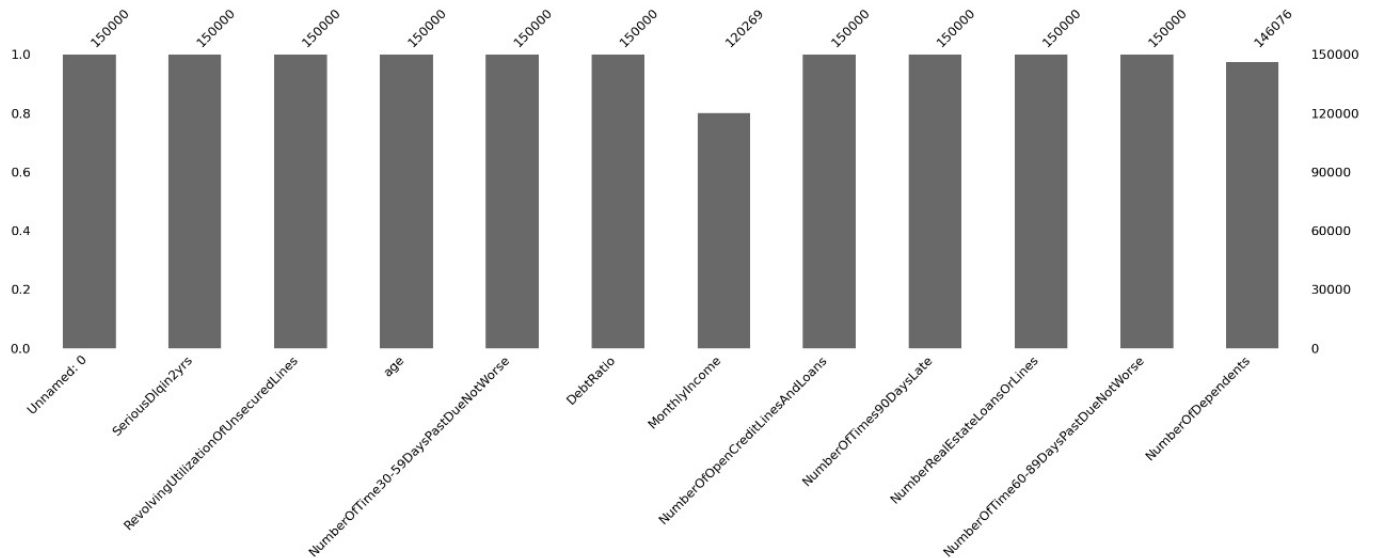


FIGURE 8 – Valeurs complètes par colonne

Pour assurer le bon apprentissage du modèle, et puisque la simple suppression des lignes avec des données manquantes entraînera à la perte de 20% de la taille initiale de la base de données, nous avons décidé de remplir ces valeurs nulles en appliquant une régression linéaire stochastique sur la colonne **NumberOfDependents** et puis sur la colonne **MonthlyIncome**. Cette régression va deviner les valeurs manquantes en se basant sur les informations apportées par les variables sans nullité. Les scores R2 de ces régressions sont respectivement 0.08 et 0.06.

Equilibrage des données d'apprentissage

Après la suppression des valeurs aberrantes, la base de données contient désormais 140.354 clients dont seulement 9025 ont **SeriousDlqin2yrs**=1. Afin d'équilibrer les données d'apprentissage, nous exécutons la requête SQL suivante :

```
SELECT *
FROM    dbo.data
WHERE   seriousdlqin2yrs = 1
UNION
SELECT TOP 9025 *
FROM    dbo.data
WHERE   seriousdlqin2yrs = 0
```

FIGURE 9 – Requête SQL de l'équilibrage

L'échantillon résultant comporte exactement le même nombre de clients ayant **SeriousDlqin2yrs**=1 que de clients avec **SeriousDlqin2yrs**=0 et a une taille globale de 18050 lignes.

Identification des meilleurs prédicteurs

On groupe la population suivant la variable **SeriousDlqin2yrs** On représente (boîtes à moustache) les deux classes pour chaque variable :

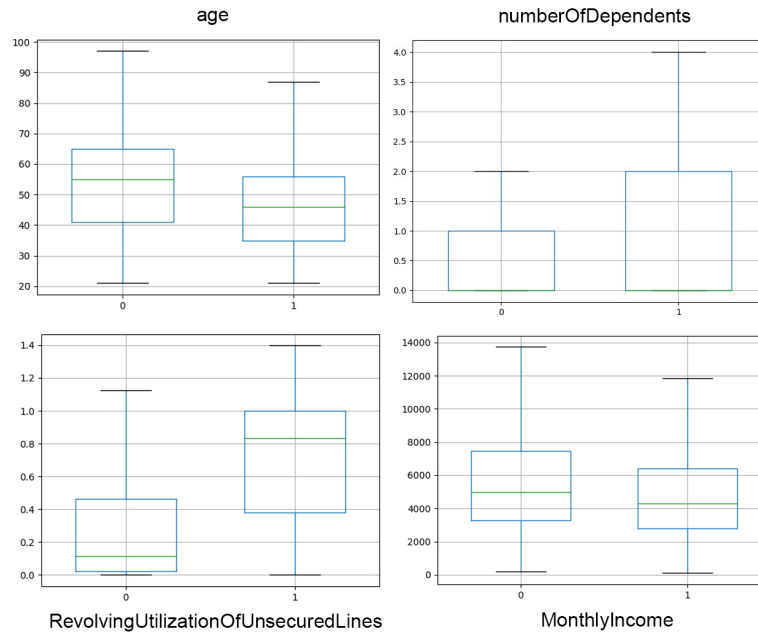


FIGURE 10 – Comparaison entre les 2 classes de clients

D'après ces représentations, nous estimons que les variables **RevolvingUtilizationOfUnsecuredLines**, **Age**, **MonthlyIncome**, **NumberOfDependents** et **NumberOfTimes90DaysLate** seront cruciales à l'identification des classes.

3.2 Modèle de prévision

Afin de valider nos modèles, nous avons choisi une méthode de validation croisée basée sur les fonctions *RepeatedStratifiedKfold* et *cross_val_score* du package **sklearnmodel_selection**. Le code est disponible en annexe.

3.2.1 Analyse factorielle discriminante

Malheureusement, nous n'avons pas trouvé d'outil prédéfini pour appliquer l'AFD, nous avons donc tenté de l'implémenter nous même. Nous avons écrit une fonction qui effectue l'apprentissage en se basant sur le modèle mathématique. La fonction calcule la barre de décision qu'on a trouvé égale à $c = -5.26e-5$.

3.2.2 Linear & Quadratic Discriminant Analysis

On applique une LDA sur les données de l'échantillon en spécifiant $Y = \text{SeriousDlqin2yrs}$ la variable à expliquer en fonction des autres. Après l'apprentissage, on choisit 300 points aléatoirement et on les visualise d'une façon à mettre en valeur la séparation des groupes.

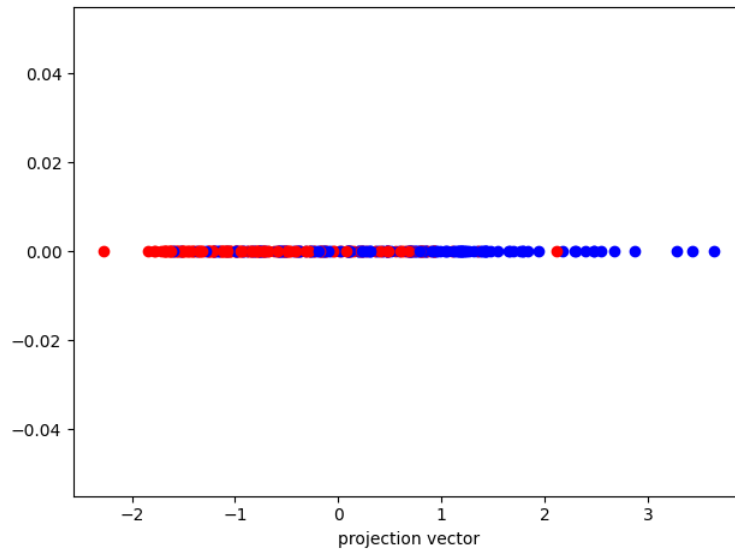


FIGURE 11 – Séparation des groupes en LDA

On effectue la QDA de la même manière.

3.2.3 Logistic Regression

On applique une régression logistique sur les données et on représente graphiquement chaque variable dans une figure.

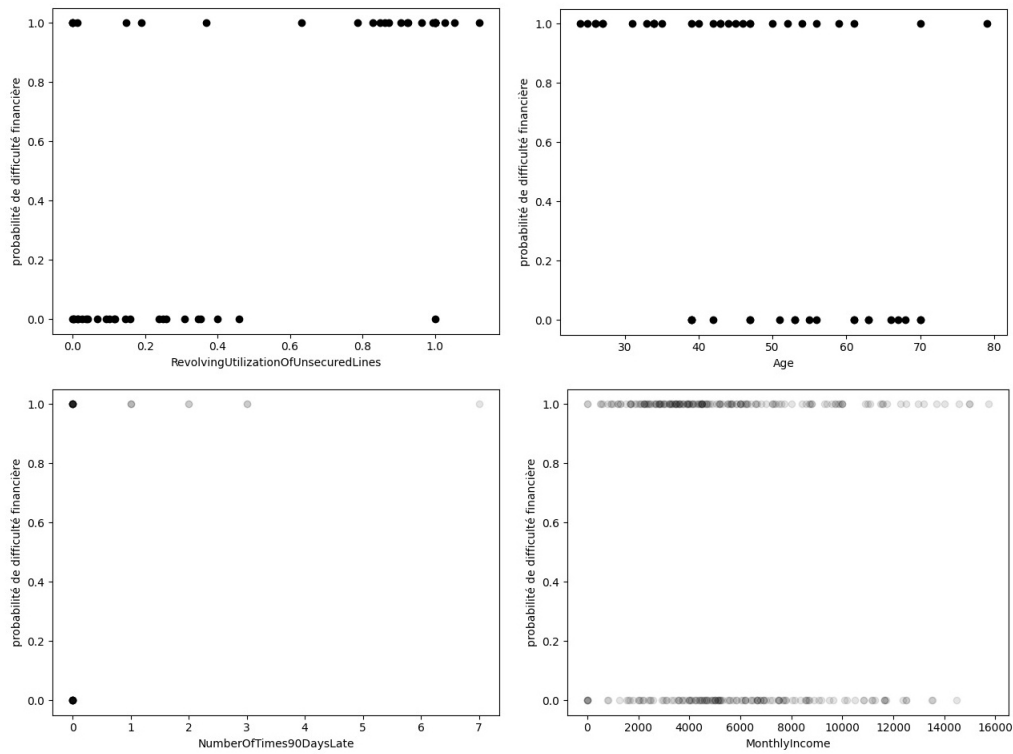


FIGURE 12 – Régression Logistique

Pour les variables **NumberOfTime90DaysLate** et **MonthlyIncome**, nous avons réduit l'opacité des points pour mieux distinguer les points superposés.

Remarquons que la variable **RevolvingUtilizationOfUnsecuredLines** prend des valeurs différentes pour chaque classe. Ceci peut être observé pour d'autres variables aussi : les personnes qui trouvent des difficultés financières (**SeriousDlqin2yrs=1**) sont généralement plus jeune et touchent en moyenne un salaire inférieur à ceux qui ne trouvent pas de problème.

Ces remarques confirment les hypothèses que nous avons fait dans la phase de prétraitement.

3.3 Phase d'évaluation et règle de décision retenue

3.3.1 Goodness of fit

On se sert de la méthode de validation définie préalablement pour calculer le score de précision de chaque modèle. Le tableau suivant décrit les résultats.

Technique	Score
LDA	0.757
QDA	0.742
Linear Regression	0.770

FIGURE 13 – Scores de précision des modèles

3.3.2 Pouvoir de prédiction

On applique nos 3 modèles sur les données du fichier *ScoringTest.csv* après l'avoir nettoyé des valeurs manquants. La fonction **confusion_matrix** permet d'obtenir les matrices de confusion suivantes :

	Prédit 0	Prédit 1
Correct 0	66666	13374
Correct 1	4	1356

FIGURE 14 – Matrice de confusion LDA

	Prédit 0	Prédit 1
Correct 0	70802	9238
Correct 1	3	1357

FIGURE 15 – Matrice de confusion QDA

	Prédit 0	Prédit 1
Correct 0	71268	8772
Correct 1	3	1357

FIGURE 16 – Matrice de confusion Logistic Regression

A partir de ces matrices de confusion on calcule le pouvoir de prédiction et on obtient les résultats suivants :

Technique	Score
LDA	0.835
QDA	0.886
Linear Regression	0.892

FIGURE 17 – Tableau des pouvoirs de prédiction

On trace également les courbes ROC de chaque modèle. Puisque la précision des modèles est très proche, les courbes ROC sont très similaires.

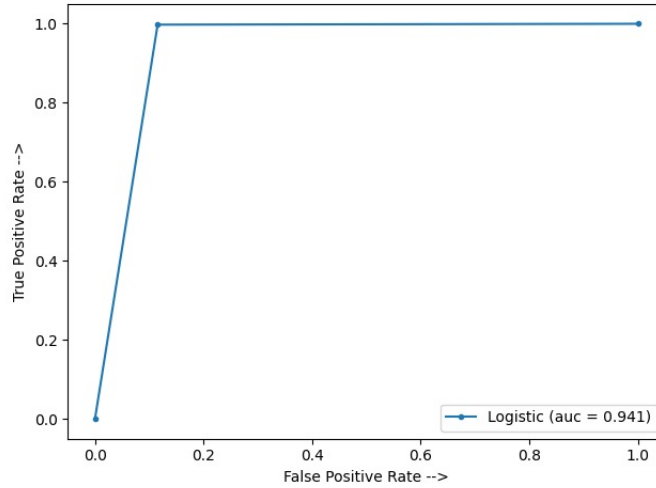


FIGURE 18 – Courbe ROC de Logistic Regression

3.3.3 Conclusion

A partir de ce résultats, on déduit que les 3 modèles donnent des résultats satisfaisants et sont capable de prédire les classes -surtout la classe 0- avec précision mais le modèle basé sur la **régression logistique** reste le plus performant.

Ce qui est logique et prévisible parce que on a les propriétés suivantes qui favorisent une régression logistique[7] :

- La variable à prédire Y est de nature dichotomique. Autrement dit, il s'agit d'un problème de classification binaire (oui/non).
- IL n'y a pas de valeurs aberrantes dans les données car elles ont été supprimé dans la phase de prétraitement.
- La corrélation entre les variables explicatives X est faible.

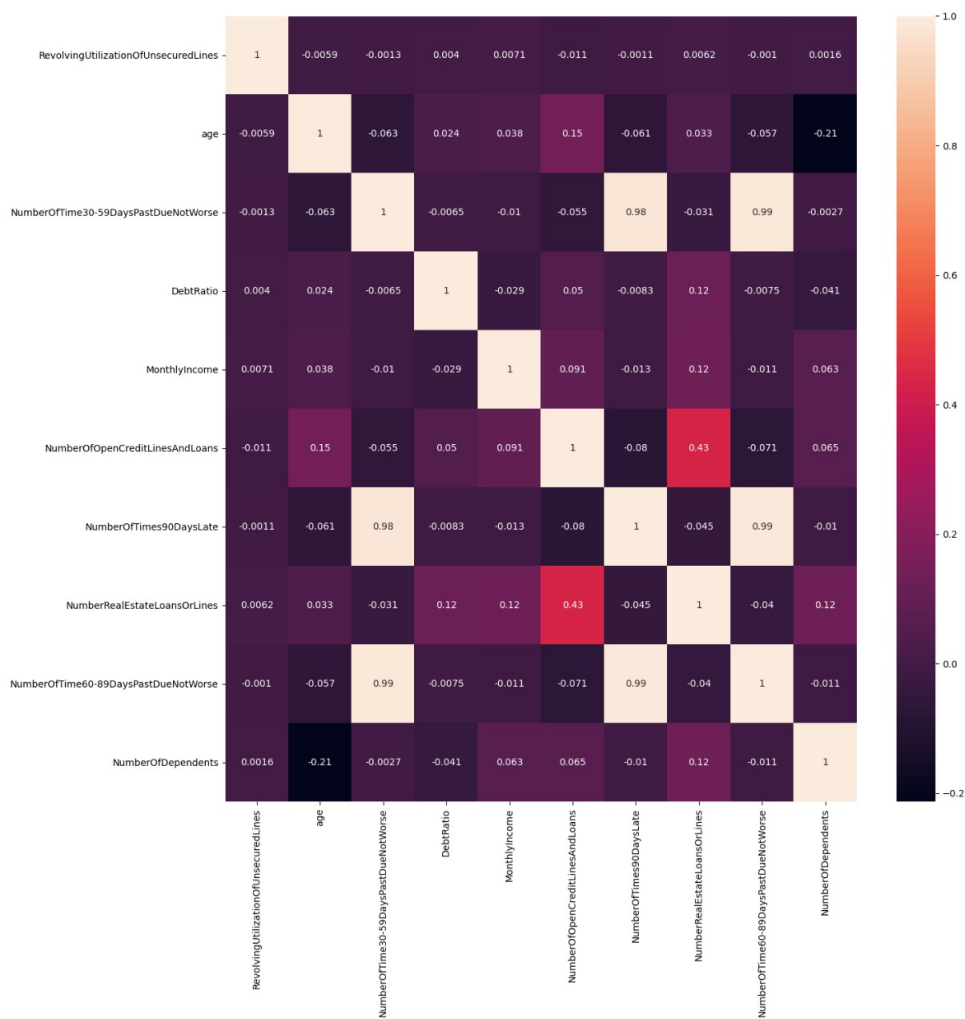


FIGURE 19 – Matrice de corrélation

4 Références bibliographiques

- [1] Missing Data Imputation using Regression - <https://www.kaggle.com/shashankasubrahmanya/missing-data-imputation-using-regression>
- [2] Pandas : Project Description - <https://pypi.org/project/pandas/>
- [3] Structed Query Language - https://fr.wikipedia.org/wiki/Structured_Query_Language
- [4] Notion De Taux De Défaut (Default Rate) - <https://aide.wesharebonds.com/hc/fr/articles/203886952-Notion-de-taux-de-d%C3%A9faut-default-rate>
- [5] Taux de défaut : Comprendre cet indicateur si complexe - <https://argent-et-salaire.com/taux-de-defaut/>
- [6] missingno - <https://github.com/ResidentMario/missingno>
- [7] What is Logistic Regression ? - <https://www.statisticssolutions.com/what-is-logistic-regression/>

5 Annexe

— Création d'une boîte à moustaches

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # changement de l'echelle de l'axe
5 plt.ylim(0, 3)
6 # lecture du fichier csv
7 data = pd.read_csv("ScoringTraining.csv")
8 # creation d'une boîte à moustache
9 boxplot = data.boxplot(column
    = "NumberOfDependents")
10 plt.show()
```

— Création de la matrice de nullité

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from quilt.data.ResidentMario import missingno_data
4 import missingno as msno
5
6 # lecture du fichier csv
7 data = pd.read_csv("ScoringTraining.csv")
8 collisions = missingno_data.nyc_collision_factors()
9 msno.matrix(data)
10 plt.show()
```

— Création des boîtes à moustaches par classe

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 |
4 data = pd.read_csv("training_data.csv")
5 # splitting dataframe by groups
6 grouped = data.groupby(data.SeriousDlqin2yrs)
7 # grouping by particular dataframe column
8 data_0 = grouped.get_group(0)
9 data_1 = grouped.get_group(1)
10
11 ax = plt.axes()
12 # first boxplot
13 boxplot = data_0.boxplot(
14     column="NumberOfDependents", positions=[1], widths=0.6, showfliers=False
15 )
16 # second boxplot
17 boxplot = data_1.boxplot(
18     column="NumberOfDependents", positions=[2], widths=0.6, showfliers=False
19 )
20 ax.set_xticklabels(["0", "1"])
21 plt.savefig("NumberOfDependents_compare.png")
```

— Linear Discriminant Analysis

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import random
4 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
5
6 # ouvrir le fichier csv
7 csv = pd.read_csv("C:\\Users\\pc\\Desktop\\training_data.csv")
8 # définition de la variable à prédire Y et des variables explicatives X
9 Y = csv.iloc[:, 1]
10 X = csv.iloc[:, 2:11]
11 # Définition du modèle
12 lda = LinearDiscriminantAnalysis(n_components=1)
13 # Apprentissage
14 score = lda.fit(X, Y).transform(X)
15 # Prédiction de Y
16 y_pred = lda.predict(X)
17 # Evaluation du modèle
18 cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
19 scores = cross_val_score(lda, X, Y, scoring="accuracy", cv=cv, n_jobs=-1)
20 print("Mean Accuracy: %.3f (%.3f)" % (numpy.mean(scores), numpy.std(scores)))
21 # Sélection aléatoire de 300 points pour la visualisation
22 random_choices = [random.randint(0, len(score) - 1) for i in range(0, 300)]
23 score1 = [score[i] for i in random_choices]
24 y_train1 = [Y[i] for i in random_choices]
25 # Représentation graphique
26 for k, i in zip(score1, y_train1):
27     if i == 1:
28         plt.scatter(k, 0, c="b")
29     else:
30         plt.scatter(k, 0, c="r")
31 plt.xlabel("projection vector")
32 plt.show()
```

— Quadratic Discrimination Analysis

```
1 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
2 import pandas as pd
3 from sklearn.model_selection import RepeatedStratifiedKFold, cross_val_score
4 # Memes étapes de LDA mais on change le type du modèle
5 csv=pd.read_csv("C:\\Users\\pc\\Desktop\\training_data.csv")
6 Y=csv.iloc[:,1]
7 X=csv.iloc[:,2:11]
8 qda=QuadraticDiscriminantAnalysis()
9 pred=qda.fit(X,Y)
10 y_pred=qda.predict(X)
11 cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
12 scores = cross_val_score(lda, X, Y, scoring='accuracy', cv=cv, n_jobs=-1)
13 print('Mean Accuracy: %.3f (%.3f)' % (numpy.mean(scores), numpy.std(scores)))
14 print("your accuracy is :"+str(accuracy))
```

— Logistic Regression

```
1 import pandas as pd
2 from sklearn.linear_model import LogisticRegression
3 import matplotlib.pyplot as plt
4 import random
5 from sklearn.metrics import accuracy_score
6 # Memes étapes de LDA mais on change le type de modèle
7 csv = pd.read_csv("C:\\Users\\pc\\Desktop\\training_data.csv")
8 Y = csv.iloc[:, 1]
9 X = csv.iloc[:, 2:11]
10 logistic_reg = LogisticRegression()
11 score = logistic_reg.fit(X, Y)
12 # On choisit une variable à visualiser
13 ruoul = csv.iloc[:, 2]
14 # Selection aléatoire de 50 points
15 random_choices=[random.randint(0,len(score)-1) for i in range(0,50)]
16 score1=[score[i] for i in random_choices]
17 y_train1=[ruoul[i] for i in random_choices]
18 # Représentation graphique
19 for k, i in zip(score, y_train1):
20     plt.scatter(i, k, c="k")
21 plt.ylabel("probabilité de difficulté financière")
22 plt.xlabel("RevolvingUtilizationOfUnsecuredLines")
23 plt.show()
```

— Prédiction des données de *ScoringTest.csv*

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import accuracy_score, r2_score, confusion_matrix
4 from sklearn.metrics import roc_curve, auc
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
7
8 csv = pd.read_csv("C:\\Users\\Achraf\\PycharmProjects\\Scoring\\training_data.csv")
9 # le fichier test.csv contient les données de ScoringTest.csv + la colonne probabilité
10 # du fichier SampleScoring.csv
11 csv2 = pd.read_csv("C:\\Users\\Achraf\\PycharmProjects\\Scoring\\test.csv")
12 Y = csv.iloc[:, 1]
13 X = csv.iloc[:, 2:11]
14 X_test = csv2.iloc[:, 2:11]
15 Y_true_tmp = csv2.iloc[:, -1]
16 # Dédduit les valeurs bool à partir des probabilités correctes
17 Y_true = []
18 for i in Y_true_tmp:
19     if i > 0.5:
20         Y_true.append(True)
21     else:
22         Y_true.append(False)
23 model = LogisticRegression()
24 # Apprentissage sur les données de training_data
25 score = model.fit(X, Y)
26 # prédiction de Y sur le fichier ScoringTest
27 Y_pred_tmp = model.decision_function(X_test)
28 Y_pred = []
29 # Dédduit les valeurs bool à partir des probabilités prédites
30 for i in Y_pred_tmp:
31     if i > 0.5:
32         Y_pred.append(True)
33     else:
34         Y_pred.append(False)
```

— Courbe ROC

```
35 logistic_fpr, logistic_tpr, threshold = roc_curve(Y_true, Y_pred)
36 auc_logistic = auc(logistic_fpr, logistic_tpr)
37 plt.figure()
38 plt.plot(
39     logistic_fpr,
40     logistic_tpr,
41     marker=".",
42     label="Logistic (auc = %0.3f)" % auc_logistic,
43 )
44 plt.xlabel("False Positive Rate -->")
45 plt.ylabel("True Positive Rate -->")
46
47 plt.legend()
48 plt.show()
```

— Matrice de Confusion

```
49 matrix = confusion_matrix(Y_true, Y_pred)
```


— Matrice de Corrélation

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sn
4
5 csv = pd.read_csv("C:\\Users\\Achraf\\PycharmProjects\\Scoring\\ScoringTraining.csv")
6 X = csv.iloc[:, 2:]
7 plt.figure(figsize=(15, 15))
8 corrMatrix = X.corr()
9 sn.heatmap(corrMatrix, annot=True)
10 plt.show()
```

— Régression linéaire Stochastique

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import missingno as msno
4 from sklearn import linear_model
5 from sklearn.svm import SVC
6 import numpy as np
7
8 data = pd.read_csv("ScoringTraining.csv")
9 colnames = list(data.columns)
10 missing_columns = ["NumberOfDependents", 'MonthlyIncome']
11
12 random_data = pd.DataFrame(columns=["Ran" + name for name in missing_columns])
13
14 for feature in missing_columns:
15     random_data["Ran" + feature] = data[feature + '_imp']
16     parameters = list(set(data.columns) - set(missing_columns) - {feature + '_imp'})
17
18     model = linear_model.LinearRegression()
19     model.fit(X=data[parameters], y=data[feature + '_imp'])
20
21     # Standard Error of the regression estimates is equal to std() of the errors of each estimates
22     predict = model.predict(data[parameters])
23     std_error = (predict[data[feature].notnull()] - data.loc[data[feature].notnull(), feature + '_imp']).std()
24
25     # observe that I preserve the index of the missing data from the original dataframe
26     random_predict = np.random.normal(size=data[feature].shape[0],
27                                       loc=predict,
28                                       scale=std_error)
29     random_data.loc[(data[feature].isnull()) & (random_predict > 0), "Ran" + feature] = random_predict[
30         (data[feature].isnull()) &
31         (random_predict > 0)]
32
33 data.to_csv(r'C:\Users\Achraf\PycharmProjects\Scoring\data_reg.csv', index = False, header=True)
```

— Analyse Factorielle Discriminante

```
1 import pandas as pd
2 import numpy as np
3 class AFD:
4     def fit(self, X, y):
5         n_features = X.shape[1] # prendre le nombre de colone existant en X
6         class_labels = np.unique(y) # vrenvoi le valeur unique du vecteur y
7
8         mean_overall = np.mean(X, axis=0) # calculer g totale
9         SW = np.zeros(
10             (n_features, n_features)
11         ) # creer une matrice de zero de n_features ligne et de n_features colone
12         SB = np.zeros((n_features, n_features))
13         gravities = []
14         for c in class_labels:
15             X_c = X[y == c] # creer la matrice x_c depuis x lorsque y=c
16             mean_c = np.mean(X_c, axis=0) # calculer la gravites des classes
17             SW = (X_c - mean_c).T.dot((X_c - mean_c)) # calculer W
18             n_c = X_c.shape[0]
19             # calculer B
20             mean_diff = mean_c - mean_overall
21             SB = n_c * (mean_diff).dot(mean_diff.T)
22             gravities.append(mean_c)
23         V = np.add(SB, SW) # calculer v
24         V_inv = np.linalg.inv(V) # calculer linv de v
25         u = np.dot(V_inv, np.subtract(gravities[0], gravities[1])) # calculer u
26         a = np.subtract(gravities[0], gravities[1]) # calculer a
27         s1 = np.dot(np.transpose(gravities[0]), np.linalg.inv(SW)) # calculer s1
28         s1 = np.dot(s1, a)
29         s2 = np.dot(np.transpose(gravities[1]), np.linalg.inv(SW)) # calculer s2
30         s2 = np.dot(s2, a)
31         c = np.add(s1, s2) * 1 / 2 # calculer c
32         print(c)
33 csv = pd.read_csv("C:\\Users\\pc\\Desktop\\training_data.csv")
34 Y = csv.iloc[:, 1]
35 X = csv.iloc[:, 2:11]
36 afd = AFD()
37 afd.fit(X, Y)
```