



## TECHNIQUES DE SCORING

**Kam Yuen CHU**  
**Ezzoubir KODAD**  
**Ikbal EL HALLAOUI**

### Responsable :

Etienne MAROT  
Pierre GEORGES

### 1. INTRODUCTION

Un score est un outil statistique et probabilistique de détection du risque. Outil d'aide à l'analyse, s'ajoutant aux renseignements sur les divers événements concernant les clients. Il est utilisé par les analystes financiers dans le diagnostic individuel de chaque client. Par ailleurs, les scores peuvent être utilisés pour apprécier le niveau de risque attaché à un ensemble de clients.

Ce projet présente la méthodologie, les résultats et les tests de la fonction score obtenue.

### 2. CONSTRUCTION DU SCORE

#### 2.1. Présentation de la méthode

La méthode de construction du score se compose de trois étapes :

- Différentes variables sont retenus en fonction de leur pouvoir discriminant,
- L'échantillon complet est séparé en un échantillon d'apprentissage, sur lequel sera fait le travail de recherche de la meilleure spécification, et en un échantillon-test qui permettra la validation du ou des modèles retenus. Les deux années 1999 et 2000 seront regroupées pour construire le fichier d'apprentissage. L'année 2001 servira d'échantillon test.
- Des classes de risque sont déterminées et des probabilités de défaillances associées sont estimées.

#### 2.2. Présentation des données et sélection des variables

La base de données correspond à des clients professionnels du Crédit Lyonnais. Elle contient 84899 individus et 18 variables pour lesquels on dispose d'information quantitatives et qualitatives, ainsi que d'un critère de défaut. Celui-ci vaut 1 quand le client est défaillant, 0 quand il est bon. Le critère de défaillance correspond au passage en contentieux dans les douze mois qui suivent le calcul du score pour éviter le problème saisonnier.

La base de données contient un ou plusieurs enregistrements pour chacun des clients (identifiés par leur IDCOM), correspondant à l'observation de leurs données sur les années 1999, 2000 et 2001. Un client présent sur les trois années aura trois

enregistrements, mais un client présent sur seulement une partie de la période peut avoir un ou deux enregistrements seulement.

Le travail préparatoire des variables s'organise en deux étapes :

- Une sélection des variables *a priori* significatives de la défaillance. Il s'agit de définir l'ensemble des variables les plus intéressantes susceptibles d'influer sur le processus de la défaillance.

Nous avons décidé de retenir toutes les variables dont on dispose<sup>1</sup>. En effet cette sélection rassemble, à la fois, des critères connus de la défaillance (nombre de jours débiteurs sur les trois derniers mois, interdiction de chéquier...), mais aussi des critères plus stratégiques pouvant avoir un rôle sur la défaillance (possession d'un compte épargne, mouvements d'affaires moyen sur les 12 derniers mois...)

- Parmi les variables précédentes, la sélection statistique identifie les plus discriminantes d'entre elles grâce à des tests du Chi-Deux<sup>2</sup> sur les quelques variables qualitatives à modalités et à des tests de variance sur les données continues de la base.

Pour des raisons de simplicité, nous avons décidé de ne pas nous lancer dans ces testes, et avons considéré qu'ils étaient vérifiés sur notre base.

### 2.3. Le traitement des valeurs extrêmes et des valeurs manquantes.

L'estimation de la fonction score nécessite naturellement de nettoyer le fichier des valeurs extrêmes et de traiter les valeurs manquantes.

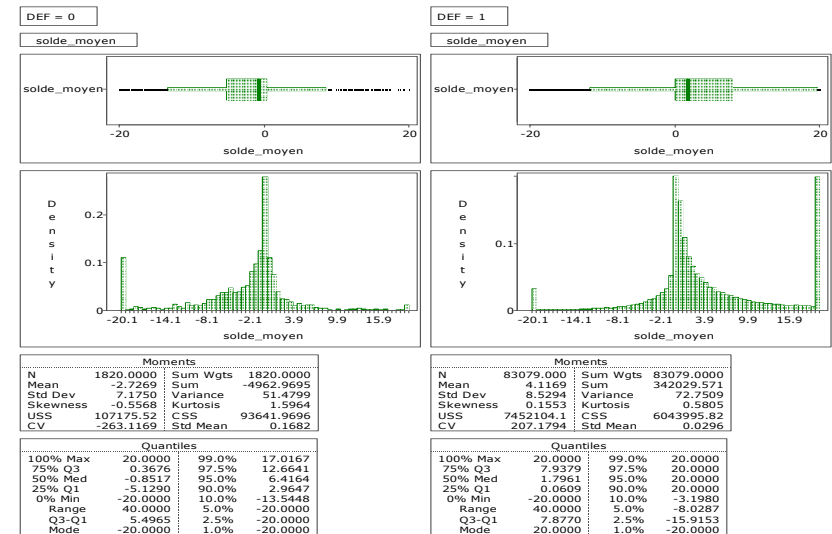
Grâce à la procédure *FREQ* du logiciel SAS nous nous sommes assurés de l'absence de valeurs manquantes. En revanche le problème du nettoyage de la base des valeurs aberrantes s'est posé. En effet, et à l'aide de la procédure *BOXPLOT*<sup>3</sup> (boîte à moustaches) de SAS nous avons pu observer la présence de valeurs extrêmes. Le détail des sorties SAS sont en annexe.

Pour illustrer cela nous présentons le cas, de la variable *solde\_moyen*.

<sup>1</sup> Bien entendu les variables *Idcom* et *Année* sont exclus de notre sélection.

<sup>2</sup> Le test du Chi 2 est un test non paramétrique permettant de tester l'hypothèse d'indépendance entre deux variables nominales

<sup>3</sup> Nous avons opté pour l'utilisation de SAS/INSIGHT.



On distingue bien les valeurs atypiques (*outliers*) situées au-delà des valeurs adjacentes (longueurs) des moustaches.

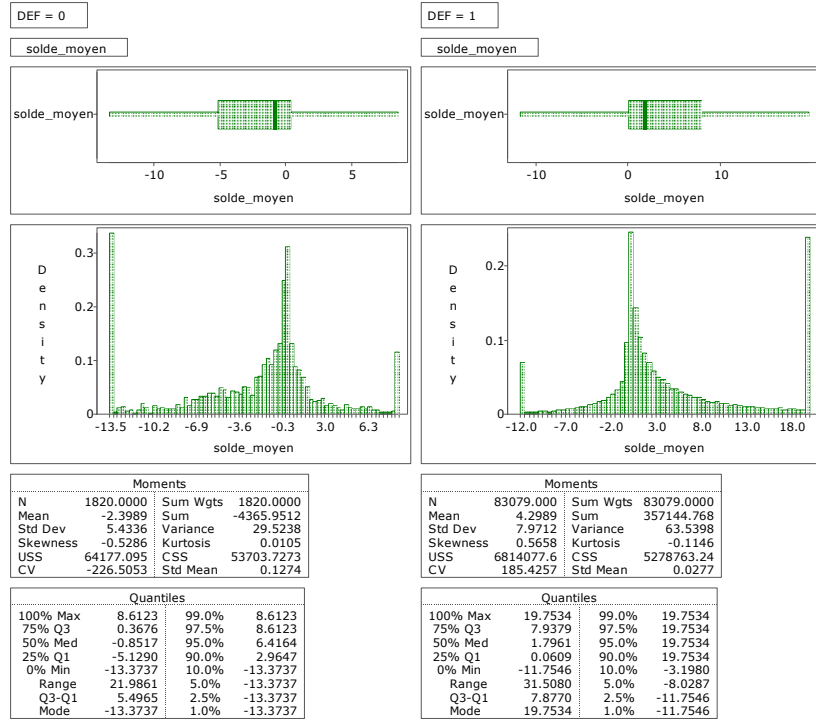
Par souci de clarté, nous allons simplement donner les équations permettant le calcul des valeurs adjacentes (borne inférieure et borne supérieure).

- Frontière basse :  $Q1 - 1,5 \cdot (Q3 - Q1)$
- Frontière haute :  $Q3 + 1,5 \cdot (Q3 - Q1)$

Généralement, lorsqu'on rencontre le problème des *outliers*, les solutions choisies sont la « *winzorisation* » et l'imputation : soit ces valeurs sont ramenées à la borne, soit elles sont « *neutralisées* » en les ramenant à la médiane. Ces solutions ne sont pas totalement satisfaisantes car l'information réelle est modifiée. Dans notre étude nous avons opté pour la première solution c'est-à-dire la *winzorisation*<sup>4</sup>.

La sortie SAS après traitement des valeurs extrêmes concernant la variable *solde\_moyen* est la suivante :

<sup>4</sup> Le programme SAS permettant de faire ces corrections est présenté en annexe.



2.4. Choix de la méthode utilisée pour construire la fonction du score.

Parmi les méthodes de scoring, la discrimination logistique à l'avantage de prendre en compte des variables aussi bien quantitatives que qualitatives. De plus elle permet d'effectuer des tests de significativité sur les paramètres et fournit aussi une aide à la sélection des variables.

Dans le cadre de cette régression, on modélise la variable DEF<sup>5</sup> qualitative (ici à deux modalités) en fonction de variables explicatives. Si p est la probabilité de la modalité {DEF=1}, on a (matriciellement) :

$$\text{Log}\left(\frac{p}{1-p}\right) = X \cdot \beta + \varepsilon$$

<sup>5</sup> A noter que nous avons effectué les changements de variable suivants : DEF=1-DEF.

Comme dans tout problème de régression, il nous faut sélectionner les variables intégrant le modèle : pour ce faire, nous disposons de deux grands types de méthode :

- Soit on les sélectionne toutes au départ, puis on retire au fur et à mesure les moins significatives (méthode dite *BACKWARD*).
- Soit on intègre les plus significatives une à une en partant du modèle ne comportant que la constante (méthode dite *FORWARD*) ; une variante consiste à ôter les variables qui ne seraient plus significatives après l'ajout de certaines autres (méthode dite *STEPWISE*).

Dans notre étude nous utiliserons la méthode *STEPWISE*.

Les deux modèles les plus classiques sont le *LOGIT*, et le *PROBIT*. Le premier suppose que les résidus suivent une loi logistique, quant au deuxième modèle la loi est supposée être une distribution normale. Cette dernière hypothèse est forte, ce qui nous entraîne à choisir le modèle *LOGIT*. Ce dernier possède théoriquement la propriété que les estimateurs des paramètres, excepté la constante, sont invariants à une surreprésentation dans l'échantillon. Ceci dit, les estimateurs de ces deux modèles sont presque identiques (proportionnels).

La sortie SAS correspondante est la suivante :

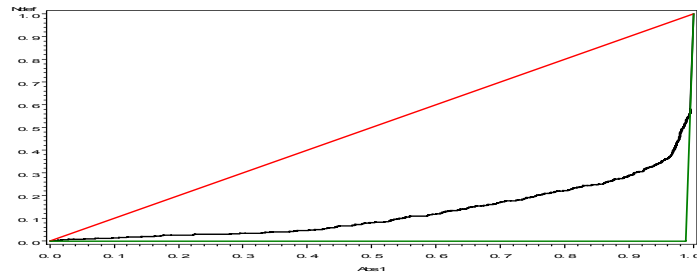
Analysis of Maximum Likelihood Estimates					
Parameter	DF	Estimate	Standard Error	wald Chi-Square	Pr > ChiSq
Intercept	1	5.0484	0.1461	1194.7896	<.0001
m1_util	1	-0.00002	6.179E-6	6.9413	0.0084
m1_util_12	1	0.000027	6.534E-6	17.3185	<.0001
ct_util	1	0.000781	0.000053	218.4857	<.0001
ct_util_12	1	-0.2574	0.0849	9.2020	0.0024
ratio1	1	-0.0121	0.00265	20.6960	<.0001
ratio3	1	0.2645	0.0149	316.3206	<.0001
anciennete	1	0.1228	0.0121	102.4686	<.0001
solde_12	1	0.0344	0.00995	11.9821	0.0005
mytaffmoy	1	0.000019	4.968E-6	15.3617	<.0001
PE	1	0.6087	0.1350	20.3198	<.0001
Pcheq	1	-2.0759	0.1051	390.2623	<.0001

Association of Predicted Probabilities and Observed Responses

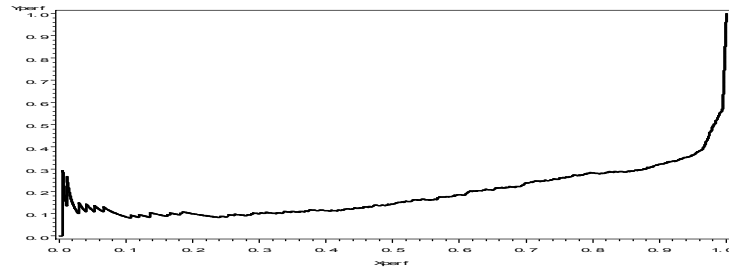
Percent Concordant	85.6	Somers' D	0.757
Percent Discordant	9.9	Gamma	0.792
Percent Tied	4.4	Tau-a	0.019
Pairs	38590050	c	0.878

Nous remarquons que la statistique Somers 'D est de 75,7% ce qui nous paraît être un bon taux, par conséquent nous avons jugé que le modèle était satisfaisant et nous n'avons pas cherché à l'améliorer en essayant de découper les variables ou de les croiser.

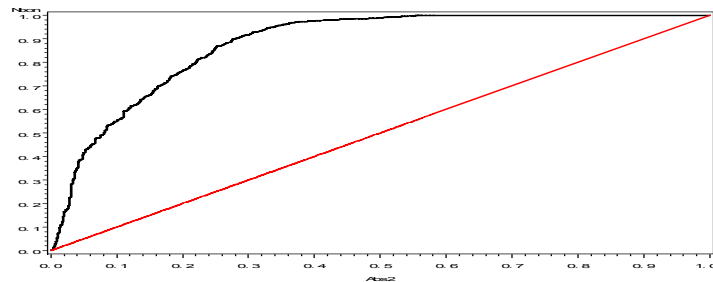
Courbe de sélection



Courbe de performance



Courbe de discrimination



3. MESURE DE LA QUALITE DE LA DISCRIMINATION

Le critère de validation de la fonction, capital dans le cas d'une fonction score, est le taux de bon classement.

Ce taux est la fréquence empirique des clients bien classés par le modèle par groupe à priori, c'est-à-dire le groupe des défaillants d'une part et le groupe des non défaillants d'autre part. un taux de bon classement global, égal au rapport du nombre de clients bien classés sur le nombre total des clients, peut être calculé, mais il peut cacher des différences entre les taux de bon classement des clients défaillants et celui des clients non défaillants ; c'est pourquoi il est préférable de le calculer sur chacun des groupes.

Les taux de bon classement reposent sur la table d'affectation suivante :

Affectation	Situation réelle	
	Défaillant	Non défaillant
Score <= seuil	n <sub>1</sub>	n <sub>3</sub>
Score > seuil	n <sub>2</sub>	n <sub>4</sub>

La règle de décision associée à la valeur du seuil permet d'affecter chaque client à un des deux groupes « défaillants » ou non défaillants. Tout client a donc un groupe auquel il appartient réellement et un groupe auquel il est affecté. Le décompte des affectations correctes, c'est-à-dire correspondant au groupe réel, fournit une estimation des taux réels de bon classement.

- Taux de bon classement réel estimé pour les clients défaillants :  $t_D = \frac{n_1}{n_1 + n_2}$
- Taux de bon classement réel estimé pour les clients non défaillants :  $t_{ND} = \frac{n_4}{n_3 + n_4}$
- Taux de bon classement réel général estimé sur l'échantillon :  $t = \frac{n_1 + n_4}{n_1 + n_2 + n_3 + n_4}$

L'étape de validation se fait à partir de l'échantillon-test.

Le tableau suivant, présente les taux de bon classement mis en évidence, d'une part, sur l'échantillon d'apprentissage constitué des données des années 1999 et 2000 et, d'autre part, sur l'échantillon de validation constitué des données 2001.

	Echantillon d'apprentissage		Echantillon test	
	Défaillant	Non défaillant	Défaillant	Non défaillant
Score <= seuil	410	1152	801	3027
Score > seuil	291	53898	318	25002

Nous avons pris comme seuil la valeur 3.

Nous en déduisons les taux de bon classement :

	Echantillon d'apprentissage	Echantillon test
$t_D$	58,49%	71,58%
$t_{ND}$	97,91%	89,20%
$t$	97,41%	88,52%

Pour que l'outil de discrimination soit efficace, il faut que, pour chacun des groupes, les taux de bon classement diffèrent sensiblement d'une répartition au hasard (taux de bon classement nettement supérieur à 50%).

Les taux de bon classement calculés ci-dessus montrent la robustesse des résultats.

#### 4. CONSTRUCTION D'UNE ECHELLE DE RATING

Pour une meilleure lecture, nous avons effectué une transformation linéaire de la fonction du score ; ceci ne change en rien sa qualité de classement.

Cette transformation<sup>6</sup> est faite de telle sorte que le score soit compris entre 0 et 1000 pour chacun des clients. La méthode de transformation est la suivante :

Soient,

$S$  : le score calculé.

$S_{\min}$  : la valeur minimale du score calculé.

$S_{\max}$  : la valeur maximale du score calculé.

$S^*$  : le score issu de la transformation.

On a, donc :  $\begin{cases} S_{\min} \leq S \leq S_{\max} \\ 0 \leq S^* \leq 1000 \end{cases}$  et pour conserver le caractère de classement de la

fonction du score, nous avons supposé une transformation linéaire, c'est-à-dire :

$S^* = a \cdot S + b$ . Il suffit donc de trouver les valeurs de  $a$  et  $b$  en résolvant le système d'équations suivant :

$$\begin{cases} S_{\min} \cdot a + b = 0 \\ S_{\max} \cdot a + b = 1000 \end{cases}$$

<sup>6</sup> Le détail du programme SAS correspondant est en annexe.

Finalement,

$$S^* = \frac{1000}{S_{\max} - S_{\min}} \cdot (S - S_{\min})$$

Nous avons ensuite découpé ce score en 10 classes de telle sorte que les probabilités de défaut par classe soient significativement différentes et croissantes.

Classe	Probabilité de défaut
950 - 1000	0,29%
894 - 950	0,86%
855 - 894	1,57%
830 - 855	2,00%
815 - 830	4,85%
800 - 815	6,28%
770 - 800	10,70%
730 - 770	13,27%
590 - 730	23,54%
0 - 590	36,66%

#### 5. CONCLUSION

L'approche probabiliste a permis de construire 10 classes de risque. Pour chacune des classes, une probabilité de défaillance est calculée, donnant ainsi une mesure de risque. Afin de valider, encore plus, la fonction de score obtenue ; il est judicieux d'explorer d'autres méthodes de scoring : la méthode des arbres, les réseaux de neurones ...

6. BIBLIORAPHIE

Pierre GEORGES, Etienne MAROT, « Techniques de scoring et applications bancaires »

Josiane CONFAIS, ISUP cours de S.A.S, « Procédures traitant les variables catégorielles : FREQ CATMOD LOGISTIC PROBIT »

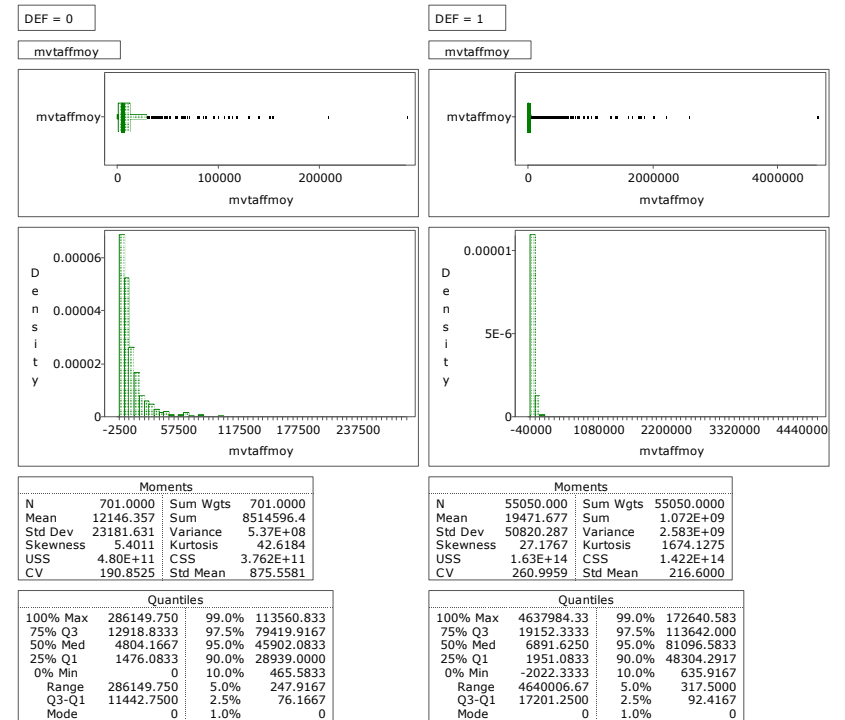
Josiane CONFAIS, ISUP cours « Variables catégorielles : modèles log-linéaire et logistique »

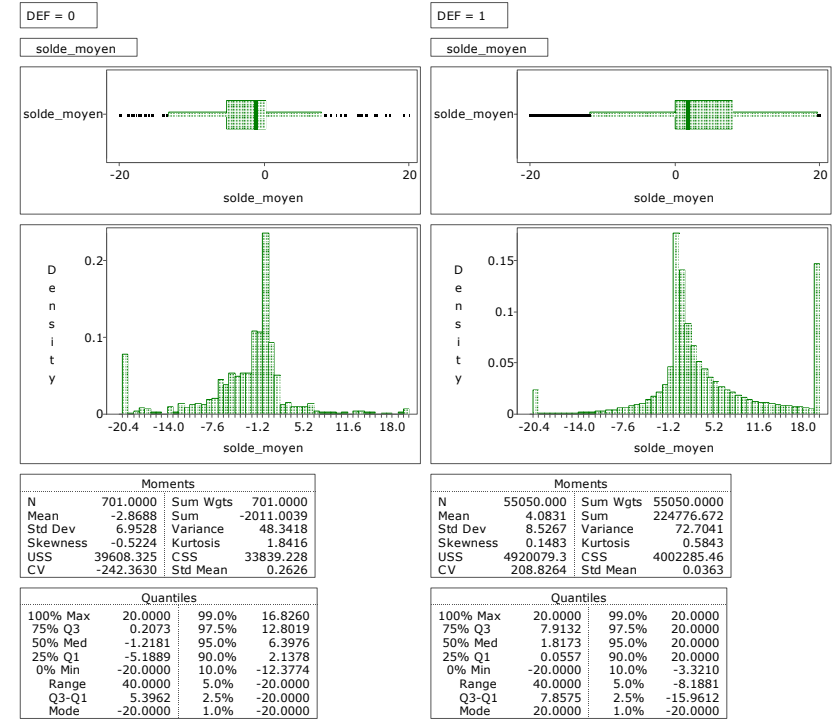
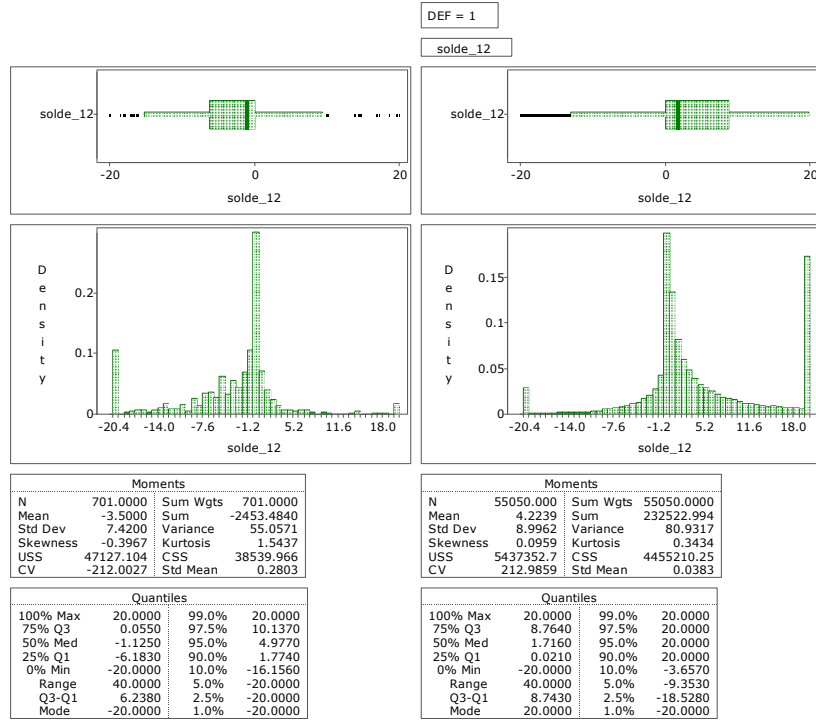
Monique LE GUEN, CNRS-MARISSE, « La boîte à moustaches de TUKEY un outil pour initier à la statistique »

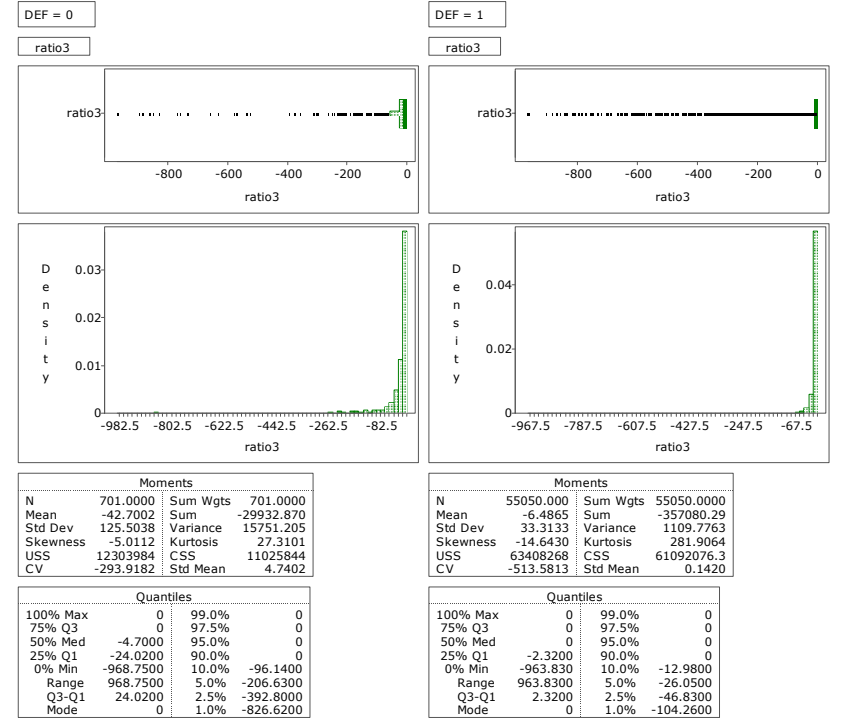
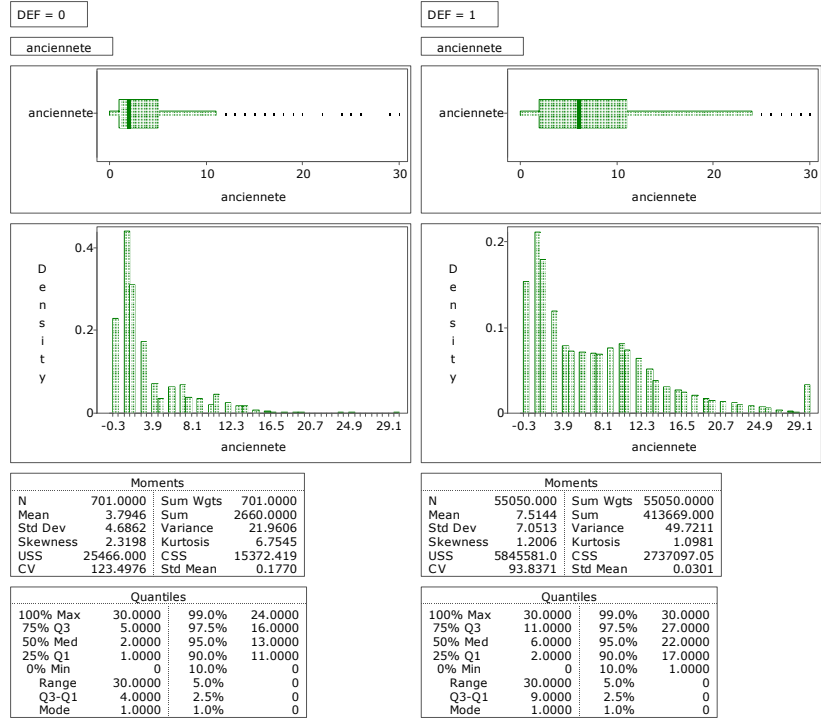
Bénédicte PALANES, BDF, « Détection précoce du risque de défaillance dans le secteur hôtels-restaurants SCORE BDFHR »

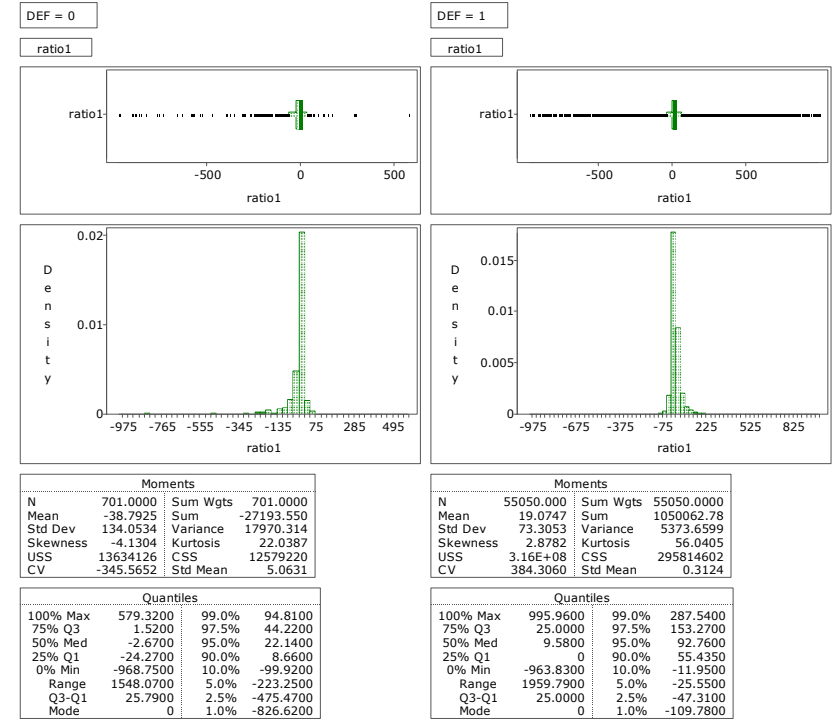
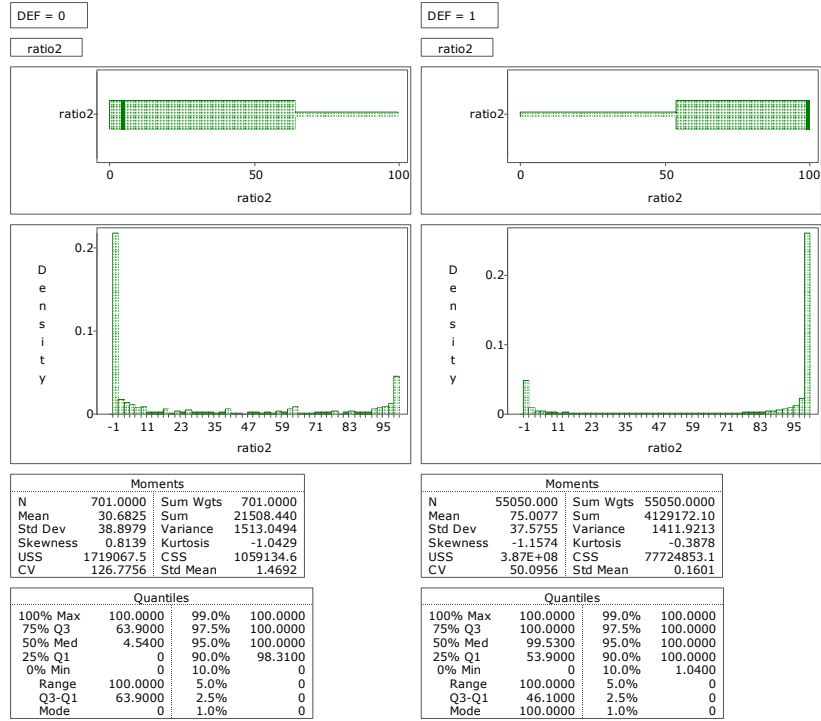
7. ANNEXE

▪ Quelques statistiques des variables explicatives.



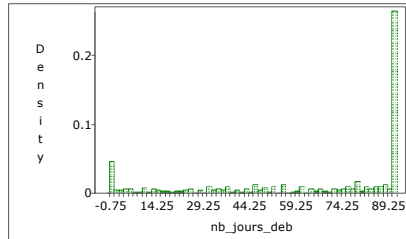
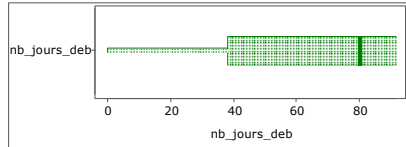






DEF = 0

nb\_jours\_deb

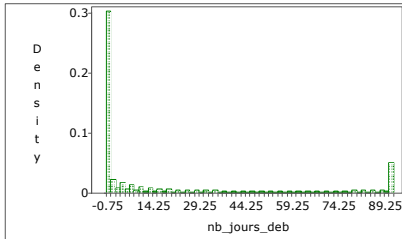
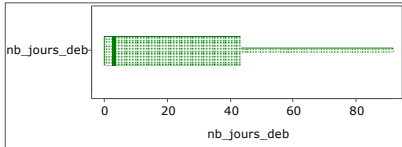


Moments			
N	701.0000	Sum Wgts	701.0000
Mean	63.7632	Sum	44698.000
Std Dev	32.7126	Variance	1070.1124
Skewness	-0.7869	Kurtosis	-0.8776
USS	3599166.0	CSS	749078.69
CV	51.3032	Std Mean	1.2355

Quantiles			
100% Max	92.0000	99.0%	92.0000
75% Q3	92.0000	97.5%	92.0000
50% Med	80.0000	95.0%	92.0000
25% Q1	38.0000	90.0%	92.0000
0% Min	0	10.0%	6.0000
Range	92.0000	5.0%	0
Q3-Q1	54.0000	2.5%	0
Mode	92.0000	1.0%	0

DEF = 1

nb\_jours\_deb

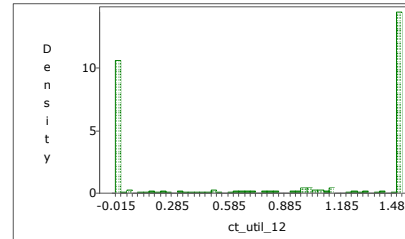
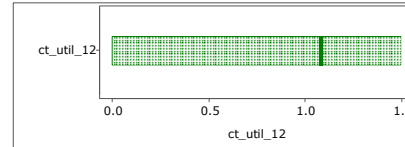


Moments			
N	55050.000	Sum Wgts	55050.0000
Mean	23.5933	Sum	1298812.00
Std Dev	32.4595	Variance	1053.6209
Skewness	1.1032	Kurtosis	-0.3484
USS	88644058	CSS	58000777.1
CV	137.5793	Std Mean	0.1383

Quantiles			
100% Max	92.0000	99.0%	92.0000
75% Q3	43.0000	97.5%	92.0000
50% Med	3.0000	95.0%	92.0000
25% Q1	0	90.0%	86.0000
0% Min	0	10.0%	0
Range	92.0000	5.0%	0
Q3-Q1	43.0000	2.5%	0
Mode	0	1.0%	0

DEF = 0

ct\_util\_12

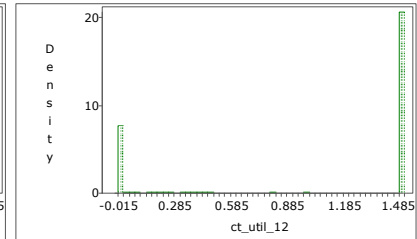
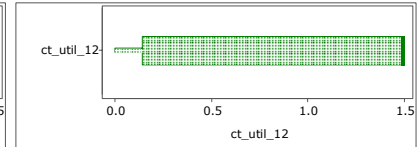


Moments			
N	701.0000	Sum Wgts	701.0000
Mean	0.8460	Sum	593.0279
Std Dev	0.6747	Variance	0.4552
Skewness	-0.2709	Kurtosis	-1.7565
USS	820.3359	CSS	318.6495
CV	79.7536	Std Mean	0.0255

Quantiles			
100% Max	1.5000	99.0%	1.5000
75% Q3	1.5000	97.5%	1.5000
50% Med	1.0805	95.0%	1.5000
25% Q1	0	90.0%	1.5000
0% Min	0	10.0%	0
Range	1.5000	5.0%	0
Q3-Q1	1.5000	2.5%	0
Mode	1.5000	1.0%	0

DEF = 1

ct\_util\_12

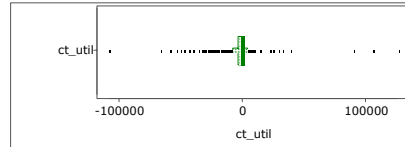


Moments			
N	55050.000	Sum Wgts	55050.0000
Mean	1.0210	Sum	56206.1692
Std Dev	0.6566	Variance	0.4312
Skewness	-0.7528	Kurtosis	-1.3126
USS	81122.414	CSS	23735.7936
CV	64.3133	Std Mean	0.0028

Quantiles			
100% Max	1.5000	99.0%	1.5000
75% Q3	1.5000	97.5%	1.5000
50% Med	1.5000	95.0%	1.5000
25% Q1	0.1394	90.0%	1.5000
0% Min	0	10.0%	0
Range	1.5000	5.0%	0
Q3-Q1	1.3606	2.5%	0
Mode	1.5000	1.0%	0

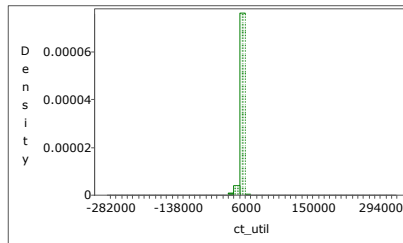
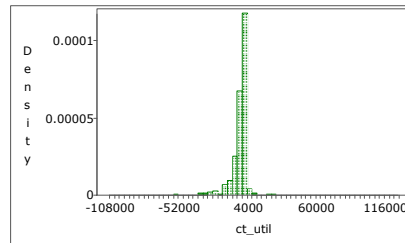
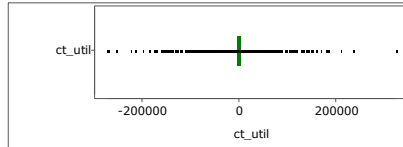
DEF = 0

ct\_util



DEF = 1

ct\_util



Moments			
N	701.0000	Sum Wgts	701.0000
Mean	-2600.343	Sum	-1822840.1
Std Dev	12045.552	Variance	1.45E+08
Skewness	1.5564	Kurtosis	45.0808
USS	1.06E+11	CSS	1.016E+11
CV	-463.2294	Std Mean	454.9542

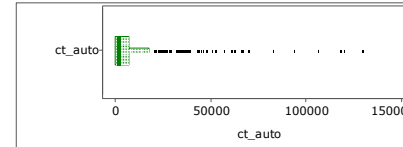
Moments			
N	55050.000	Sum Wgts	55050.0000
Mean	-1307.540	Sum	-71980052
Std Dev	8749.1990	Variance	76548483.8
Skewness	-2.5640	Kurtosis	210.1465
USS	4.21E+12	CSS	4.214E+12
CV	-669.1346	Std Mean	37.2898

Quantiles			
100% Max	127447.333	99.0%	25442.4167
75% Q3	0	97.5%	7418.5000
50% Med	0	95.0%	1610.5833
25% Q1	-3260.5833	90.0%	0
0% Min	-107592.17	10.0%	-9564.2857
Range	235039.500	5.0%	-16562.333
Q3-Q1	3260.5833	2.5%	-29109.250
Mode	0	1.0%	-47304.167

Quantiles			
100% Max	325186.583	99.0%	9193.3333
75% Q3	0	97.5%	1174.0000
50% Med	0	95.0%	0
25% Q1	-195.2500	90.0%	0
0% Min	-271298.00	10.0%	-3848.0833
Range	596484.583	5.0%	-8615.0833
Q3-Q1	195.2500	2.5%	-16204.750
Mode	0	1.0%	-30190.750

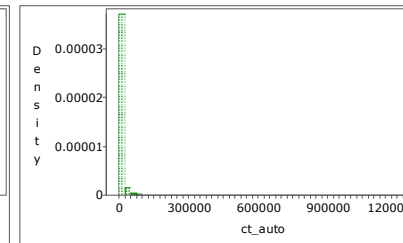
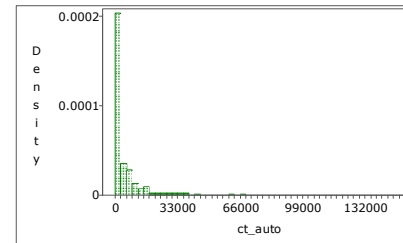
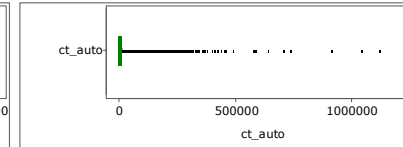
DEF = 0

ct\_auto



DEF = 1

ct\_auto



Moments			
N	701.0000	Sum Wgts	701.0000
Mean	7152.5813	Sum	5013959.5
Std Dev	16053.951	Variance	2.58E+08
Skewness	4.5320	Kurtosis	26.4098
USS	2.16E+11	CSS	1.804E+11
CV	224.4497	Std Mean	606.3494

Moments			
N	55050.000	Sum Wgts	55050.0000
Mean	6749.1748	Sum	371542075
Std Dev	23880.041	Variance	570256374
Skewness	14.4009	Kurtosis	432.1434
USS	3.39E+13	CSS	3.139E+13
CV	353.8216	Std Mean	101.7786

Quantiles			
100% Max	151895.250	99.0%	82831.0000
75% Q3	7241.0000	97.5%	52595.2500
50% Med	1524.0000	95.0%	34936.2500
25% Q1	0	90.0%	20326.3333
0% Min	0	10.0%	0
Range	151895.250	5.0%	0
Q3-Q1	7241.0000	2.5%	0
Mode	0	1.0%	0

Quantiles			
100% Max	1242569.00	99.0%	93375.5833
75% Q3	4573.0000	97.5%	57897.0833
50% Med	0	95.0%	33539.0000
25% Q1	0	90.0%	15727.3333
0% Min	0	10.0%	0
Range	1242569.00	5.0%	0
Q3-Q1	4573.0000	2.5%	0
Mode	0	1.0%	0

