

MTH8302
Modèles de régression et d'analyse de la variance
Devoir 2
distribution : 31 mai 2018
remise : 12 juin 2018 - 23h59 (au plus tard)

Ce travail est réalisé individuellement par chaque étudiant inscrit au cours.

Chaque étudiant le fait **SEUL** sans demander de l'aide à d'autres.

En apposant sa signature ci-dessous, l'étudiant (e) certifie sur son honneur avoir fait ce travail **SEUL**.

L'obtention des résultats présentés et la rédaction de ce travail ne fait l'objet d'aucun plagiat, partiel ou total.

Information concernant le plagiat à Polytechnique : <http://www.polymtl.ca/etudes/ppp/index.php>

Exigences pour la rédaction du rapport consulter la page 4 du plan de cours

<http://www.groupe.polymtl.ca/mth6301/mth8302/Autres/2018-MTH8302-PlanCours.pdf>

Compléter l'information suivante et **transmettez cette page comme la page 1** de votre rapport de devoir.

MTH8302 Modèles de régression et d'analyse de variance

NOM BETTACHE PRÉNOM Lyes Heythem

MATRICULE 1923715 SIGNATURE 

- Transmettre votre rapport par courriel à bernard.clement@polymtl.ca
- Nom suggéré pour le fichier à transmettre : NomFamille-matricule-MTH8302-Devoir2.pdf

TABLEAU CORRECTION

| | valeur | obtenu |
|---------------------------|------------|--------|
| No 5-BostonHousing | 30 | |
| No 6-BodyFat-Femme | 30 | |
| No 7-Penta | 30 | |
| Qualité | 10 | |
| TOTAL | 100 | |

- Les données pour la réalisation du devoir sont disponibles sur le site WEB du cours

<http://www.groupe.polymtl.ca/mth6301/MTH8302.htm/>

Remarque : dans la 1^{er} partie excersice 1, quand on a utilisé statistica pour développer la méthode *Forward Stepwise (ou Backward Stepwise)*, il y a 2 chemin :
 1-Advanced Models->General regression->Multiple regression->Option/Forward(Backward)
 2- Multiple regression->Advanced option-> Forward(Backward)
 J'ai utilisé les 2 et j'ai trouvé des résultats différent (dans ce rapport j'ai utilisé chemin 2)

No 5 Étude de modélisation avec plusieurs méthodes

Données = BostonHousing.sta

Réponse

5a) Modèle de Régression Ordinaire (MRO) (RÉGRESSION LINÉAIRE MULTIPLE)

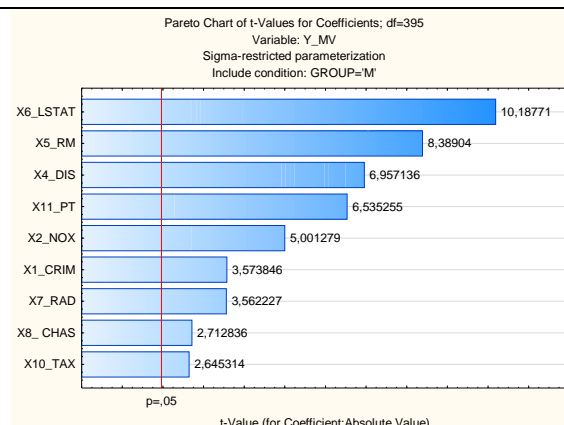
$$Y = b_0 + b_1x_1 + b_2x_2 + \dots + b_{11}x_{11} + e$$

| Var | Nom | coefficient | MRO ordinaire |
|-----|--------------------|----------------------------------|---------------|
| X0 | GENERAL intercepte | b0 | 41,4087 |
| X1 | CRIM | b1 | -0,1251 |
| X2 | NOX | b2 | -21,1251 |
| X3 | AGE | b3 | 0,0040 |
| X4 | DIS | b4 | -1,2916 |
| X5 | RM | b5 | 3,9002 |
| X6 | LSTAT | b6 | -0,5639 |
| X7 | RAD | b7 | 0,2685 |
| X8 | CHAS | b8 | 2,7103 |
| X9 | INDUS | b9 | -0,0022 |
| X10 | TAX | b10 | -0,0102 |
| X11 | PT | b11 | -0,9708 |
| | | SS resid résiduelle | 9826,13 |
| | | MSE = sigma ² (ANOVA) | 25,003 |
| | | R ² | 0,71963493 |
| | | R ² ajusté | 0,71178756 |

Ce tableau contient les coefficients de notre modèle Y basé sur les 11 variables X1,..., X11

| Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in 2018-MTH8302-Devoirs-Data) | | | | | | |
|--|-----------|-----------------|----------|----------------|----------|----------|
| R= .84831299 R ² = .71963493 Adjusted R ² = .71178756 | | | | | | |
| F(11,393)=91.704 p<0.0000 Std Error of estimate: 5.0003 | | | | | | |
| Include condition: GROUP=M' | | | | | | |
| N=405 | b* | Std. Err. of b* | b | Std. Err. of b | t(393) | p-value |
| Intercept | | | 41.4087 | 5.811334 | 7.12551 | 0.000000 |
| X1_CRIM | -0.123260 | 0.034611 | -0.1251 | 0.035134 | -3.56129 | 0.000414 |
| X2_NOX | -0.266628 | 0.056691 | -21.1251 | 4.491677 | -4.70317 | 0.000004 |
| X3_AGE | 0.012167 | 0.044953 | 0.0040 | 0.014906 | 0.27067 | 0.786790 |
| X4_DIS | -0.294117 | 0.047968 | -1.2916 | 0.210649 | -6.13146 | 0.000000 |
| X5_RM | 0.296059 | 0.036317 | 3.9002 | 0.480050 | 8.12453 | 0.000000 |
| X6_LSTAT | -0.432853 | 0.044747 | -0.5639 | 0.058293 | -9.67325 | 0.000000 |
| X7_RAD | 0.249104 | 0.072119 | 0.2685 | 0.077723 | 3.45408 | 0.000612 |
| X8_CHAS | 0.073912 | 0.027634 | 2.7103 | 1.013346 | 2.67462 | 0.007794 |
| X9_INDUS | -0.001657 | 0.051909 | -0.0022 | 0.070207 | -0.03192 | 0.974552 |
| X10_TAX | -0.183696 | 0.076450 | -0.0102 | 0.004236 | -2.40284 | 0.016731 |
| X11_PT | -0.223687 | 0.034786 | -0.9708 | 0.150978 | -6.43035 | 0.000000 |

D'après les P-value on remarque que tous les variables Xi sont significatives sauf les variables X3_AGE et X9 INDUS qui ne sont pas Significatives. (Variables Significatives p-level≤ 0,05)



D'après ce graphe on a confirmé la remarque de tableau précédant

On remarque que R² élevé (0.7196) et R²aj légèrement inférieur à R²

5b)

Les données ne présentent pas un problème de multi colinéarité

La preuve :

Pour trouver si nous avons la multicollinéarité on a utilisé les critères de Détection multicollinéarité

Critère 1 : matrice de corrélation des variables X

| Correlations (BostonHousing.sta in 2018-MTH8302-Devoirs-Data) | | | | | | | | | | | | |
|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| Marked correlations are significant at $p < .05000$ | | | | | | | | | | | | |
| N=405 (Casewise deletion of missing data) | | | | | | | | | | | | |
| Include condition: GROUP=M | | | | | | | | | | | | |
| | X1_CRIM | X2_NOX | X3_AGE | X4_DIS | X5_RM | X6_LSTAT | X7_RAD | X8_CHAS | X9_NDUS | X10_TAX | X11_PT | Y_MV |
| Variable | | | | | | | | | | | | |
| X1_CRIM | 1.000000 | 0.405092 | 0.341859 | -0.371778 | -0.183302 | 0.424085 | 0.608001 | -0.049381 | 0.394384 | 0.566047 | 0.289547 | -0.377012 |
| X2_NOX | 0.405092 | 1.000000 | 0.721333 | -0.770838 | -0.289010 | 0.590483 | 0.610317 | 0.117650 | 0.761074 | 0.666941 | 0.164000 | -0.421667 |
| X3_AGE | 0.341859 | 0.721333 | 1.000000 | -0.728034 | -0.231930 | 0.596188 | 0.458214 | 0.111474 | 0.633658 | 0.501316 | 0.243265 | -0.359839 |
| X4_DIS | -0.371778 | -0.770838 | -0.728034 | 1.000000 | 0.192688 | -0.484251 | -0.502597 | -0.117004 | -0.705015 | -0.540198 | -0.220393 | 0.230695 |
| X5_RM | -0.183302 | -0.289010 | -0.231930 | 0.192688 | 1.000000 | -0.605808 | -0.174608 | 0.077613 | -0.386300 | -0.265062 | -0.341690 | 0.685445 |
| X6_LSTAT | 0.424085 | 0.590483 | 0.596188 | -0.484251 | -0.605808 | 1.000000 | 0.469223 | -0.042379 | 0.596729 | 0.528165 | 0.377576 | -0.740351 |
| X7_RAD | 0.608001 | 0.610317 | 0.458214 | -0.502597 | -0.174608 | 0.469223 | 1.000000 | 0.006493 | 0.595796 | 0.909233 | 0.480713 | -0.364852 |
| X8_CHAS | -0.049381 | 0.117650 | 0.111474 | -0.117004 | 0.077613 | -0.042379 | 0.006493 | 1.000000 | 0.081321 | -0.013160 | -0.109490 | 0.154035 |
| X9_NDUS | 0.394384 | 0.761074 | 0.633658 | -0.705015 | -0.386300 | 0.596729 | 0.595796 | 0.081321 | 1.000000 | 0.715589 | 0.364074 | -0.468867 |
| X10_TAX | 0.566047 | 0.666941 | 0.501316 | -0.540198 | -0.265062 | 0.528165 | 0.909233 | -0.013160 | 0.715589 | 1.000000 | 0.462376 | -0.452231 |
| X11_PT | 0.289547 | 0.164000 | 0.243265 | -0.220393 | -0.341690 | 0.377576 | 0.480713 | -0.109490 | 0.364074 | 0.462376 | 1.000000 | -0.473461 |
| Y_MV | -0.377012 | -0.421667 | 0.359839 | 0.230695 | 0.685445 | -0.740351 | -0.364852 | 0.154035 | -0.468867 | -0.452231 | -0.473461 | 1.000000 |

D'après le matrice de corrilation les données ne présentent pas le problème de multi colinéarité, si on prendre $rij \geq 0.95$. Puisque le critère 1 pour la détection de multicolonéarité n'est pas presente ($rij < 0.95$)

($R = (rij)$ matrice de corrélation des variables X critère1 $ri \geq 0.95$ nécessaire mais non suffisant) Mais on remarque pour $rij \geq 0.70$ il y a une corrélation entre quelque variable donc, on vérifie les deux autres critères

Critère 2 : Variance inflation factors

| Collinearity statistics for terms in the equation (BostonHousing.sta in 2018-MTH8302-Devoirs-Data) | | | | | | | | | |
|--|-----------|--------------------|----------|--------------|--------------|---------------|-----------|----------|--|
| Sigma-restricted parameterization | | | | | | | | | |
| Include condition: GROUP=M | | | | | | | | | |
| Effect | Tolerance | Variance Inflation | R square | Y_MV Beta in | Y_MV Partial | Y_MV Semi-par | Y_MV t | Y_MV p | |
| X1_CRIM | 0.595250 | 1.679150 | 0.404475 | -0.123290 | -0.176813 | -0.095120 | -3.561295 | 0.000414 | |
| X2_NOX | 0.221972 | 4.504568 | 0.778072 | -0.766824 | -0.230836 | -0.125513 | -4.703169 | 0.000036 | |
| X3_AGE | 0.353038 | 2.832589 | 0.646962 | 0.012167 | 0.013620 | 0.007224 | 0.270666 | 0.786786 | |
| X4_DIS | 0.310042 | 3.225367 | 0.689577 | -0.294115 | -0.295481 | -0.163768 | -5.131464 | 0.000000 | |
| X5_RM | 0.540826 | 1.848758 | 0.459107 | 0.295056 | 0.379217 | 0.217002 | 8.124529 | 0.000000 | |
| X6_LSTAT | 0.365265 | 2.806717 | 0.643715 | -0.432633 | -0.436296 | -0.258367 | -9.673248 | 0.000000 | |
| X7_RAD | 0.137152 | 7.298510 | 0.862873 | 0.249104 | 0.171649 | 0.092269 | 3.454085 | 0.000513 | |
| X8_CHAS | 0.934178 | 1.074588 | 0.065820 | 0.073911 | 0.133705 | 0.071437 | 2.674620 | 0.007798 | |
| X9_NDUS | 0.264751 | 3.777131 | 0.735248 | -0.001659 | -0.001610 | -0.000826 | -0.031920 | 0.974552 | |
| X10_TAX | 0.122062 | 8.192532 | 0.877937 | -0.183698 | -0.120326 | -0.064178 | -2.402838 | 0.016731 | |
| X11_PT | 0.589547 | 1.692162 | 0.410452 | -0.223687 | -0.308542 | -0.171754 | -6.30351 | 0.000000 | |

D'après le tableau on a max VIF =9,1925<10 donc les données ne présentent pas un problème de multicollinéarité (max VIF ≥ 10 c-à-d $R^2 \geq 0,90$)

Critère 3 : Indice Conditionnement

| Value number | Eigenvalue | % Total variance | Cumulative Eigenvalue | Cumulative % | IC |
|--------------|------------|------------------|-----------------------|--------------|------------|
| 1 | 5,462799 | 49,66181 | 5,46280 | 49,6618 | 1 |
| 2 | 1,372640 | 12,47854 | 6,83544 | 62,1404 | 3,97977587 |
| 3 | 1,129551 | 10,26864 | 7,96499 | 72,4090 | 4,83625771 |
| 4 | 0,845876 | 7,68978 | 8,81087 | 80,0988 | 6,45815553 |
| 5 | 0,655850 | 5,96228 | 9,46672 | 86,0611 | 8,32933748 |
| 6 | 0,500154 | 4,54685 | 9,96687 | 90,6079 | 10,9222384 |
| 7 | 0,348620 | 3,16927 | 10,31549 | 93,7772 | 15,669793 |
| 8 | 0,240779 | 2,18890 | 10,55627 | 95,9661 | 22,6880139 |
| 9 | 0,211233 | 1,92030 | 10,76750 | 97,8864 | 25,8615174 |
| 10 | 0,163825 | 1,48932 | 10,93133 | 99,3757 | 33,3452692 |
| 11 | 0,068673 | 0,62430 | 11,00000 | 100,0000 | 79,5475137 |

On a vérifié aussi le Critère 3 ou on a trouvé que les données ne présentent pas un problème de multicollinéarité ($IC = \lambda_1/\lambda_k < 100$ $k = 2, 3, \dots$)

D'après les 3 critères les données ne présentent pas le problème de multi colinéarité.

5c)

Modèle de Régression avec Sélection pas à pas Avant (Forward Stepwise) (MRF)

$$Y = b_0 + b_1 \cdot x_1 + b_2 \cdot x_2 + b_4 \cdot x_4 \dots + b_8 \cdot x_8 + b_{10} \cdot x_{10} + b_{11} \cdot x_{11} + e$$

Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
R= .84828194 R²= .71958226 Adjusted R²= .71319299
F(9,395)=112.62 p<0.0000 Std Error of estimate: 4.9881
Include condition: GROUP=M'

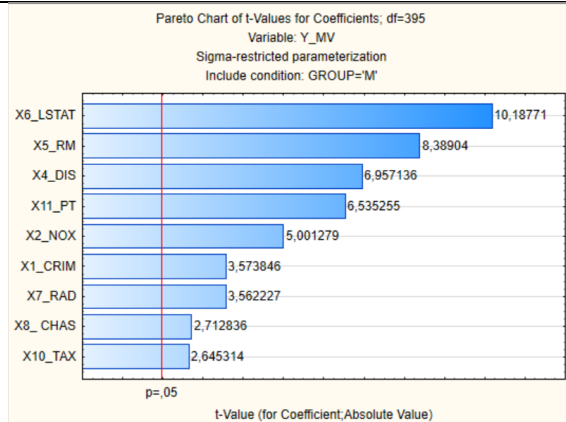
| | b* | Std Err. of b* | b | Std Err. of b | t(395) | p-value |
|-----------|-----------|----------------|----------|---------------|----------|----------|
| Intercept | | | 41.3312 | 5.781619 | 7.1487 | 0.000000 |
| X6_LSTAT | -0.428992 | 0.042109 | -0.5589 | 0.054856 | -10.1877 | 0.000000 |
| X5_RM | 0.296865 | 0.035387 | 3.9241 | 0.467760 | 8.3890 | 0.000000 |
| X11_PT | -0.222863 | 0.034102 | -0.9673 | 0.148007 | -6.5353 | 0.000000 |
| X4_DIS | -0.298273 | 0.042873 | -1.3098 | 0.198273 | -6.9571 | 0.000000 |
| X2_NOX | -0.263239 | 0.052634 | -20.8566 | 4.170248 | -5.0013 | 0.000001 |
| X8_CHAS | 0.074347 | 0.027406 | 2.7263 | 1.004961 | 2.7128 | 0.006963 |
| X1_CRIM | -0.123308 | 0.034503 | -0.1252 | 0.035024 | -3.5738 | 0.000395 |
| X7_RAD | 0.248692 | 0.069814 | 0.2680 | 0.075239 | 3.5622 | 0.000413 |
| X10_TAX | -0.184822 | 0.069868 | -0.102 | 0.003871 | -2.6453 | 0.008487 |

Summary of stepwise regression, variable: Y_MV (BostonHousing.sta in 2018-MTH8302-Devoirs-Data)
Forward stepwise
P to enter: .05, P to remove: .05
Include condition: GROUP=M'

| Effect | Steps | Degr. of Freedom | F to remove | P to remove | F to enter | P to enter | Effect status |
|----------|---------------|------------------|-------------|-------------|------------|------------|---------------|
| X1_CRIM | Step Number 1 | 1 | | | 66.7724 | 0.000000 | Out |
| X2_NOX | | 1 | | | 87.1500 | 0.000000 | Out |
| X3_AGE | | 1 | | | 59.9438 | 0.000000 | Out |
| X4_DIS | | 1 | | | 22.6534 | 0.000003 | Out |
| X5_RM | | 1 | | | 357.1414 | 0.000000 | Out |
| X6_LSTAT | | 1 | | | 488.8275 | 0.000000 | Entered |
| X7_RAD | | 1 | | | 61.8841 | 0.000000 | Out |
| X8_CHAS | | 1 | | | 9.7942 | 0.001878 | Out |
| X9 INDUS | | 1 | | | 113.5585 | 0.000000 | Out |
| X10_TAX | | 1 | | | 103.6076 | 0.000000 | Out |
| X11_PT | | 1 | | | 116.4407 | 0.000000 | Out |

...

| Effect | Steps | Degr. of Freedom | F to remove | P to remove | F to enter | P to enter | Effect status |
|----------|----------------|------------------|-------------|-------------|------------|------------|---------------|
| X3_AGE | | 1 | | | 0.0833 | 0.773004 | Out |
| X6_LSTAT | Step Number 10 | 1 | 103.7894 | 0.000000 | | | In |
| X5_RM | | 1 | 70.3760 | 0.000000 | | | In |
| X11_PT | | 1 | 42.7096 | 0.000000 | | | In |
| X4_DIS | | 1 | 48.4017 | 0.000000 | | | In |
| X2_NOX | | 1 | 25.0128 | 0.000001 | | | In |
| X8_CHAS | | 1 | 7.3595 | 0.006963 | | | In |
| X1_CRIM | | 1 | 12.7724 | 0.000395 | | | In |
| X7_RAD | | 1 | 12.6895 | 0.000413 | | | In |
| X10_TAX | | 1 | 6.9977 | 0.008487 | | | In |
| X9 INDUS | | 1 | | | 0.0006 | 0.980928 | Out |
| X3_AGE | | 1 | | | 0.0730 | 0.787163 | Out |



On remarque que la méthode de Régression avec Forward Stepwise a éliminé (négligé) les variables X3_AGE et X9 INDUS qui ne sont pas significatives.

Forward Stepwise ne vérifie pas l'états de toutes les variables après la sélection chaque step.

d)

Modèle de Régression avec Sélection pas à pas Arrière (Backward Stepwise) (MRB)

Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
R= .83553061 R²= .69811140 Adjusted R²= .69432833
F(5,399)=184.54 p<0.0000 Std Error of estimate: 5.1495
Include condition: GROUP=M'

| | b* | Std Err. of b* | b | Std Err. of b | t(399) | p-value |
|-----------|-----------|----------------|----------|---------------|----------|----------|
| Intercept | | | 36.8479 | 5.268551 | 6.9939 | 0.000000 |
| X2_NOX | -0.246248 | 0.047446 | -19.5103 | 3.759152 | -5.1991 | 0.000000 |
| X4_DIS | -0.291856 | 0.043995 | -1.2817 | 0.193201 | -6.6336 | 0.000000 |
| X5_RM | 0.316397 | 0.035411 | 4.1822 | 0.468074 | 8.9350 | 0.000000 |
| X6_LSTAT | -0.463725 | 0.042430 | -0.6041 | 0.055274 | -10.9291 | 0.000000 |
| X11_PT | -0.214198 | 0.030493 | -0.9297 | 0.132344 | -7.0246 | 0.000000 |

Analysis of Variance; DV: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
Include condition: GROUP=M'

| Effect | Sums of Squares | df | Mean Squares | F | p-value |
|----------|-----------------|-----|--------------|----------|---------|
| Regress. | 24467.15 | 5 | 4893.429 | 184.5359 | 0.00 |
| Residual | 10580.46 | 399 | 26.517 | | |
| Total | 35047.62 | | | | |

On remarque que la méthode de Régression avec Backward Stepwise a éliminé les variables X1_CRIM, X3_AGE, X7_RAD, X8_CHAS, X9 INDUS et X10_TAX qui ne sont pas significatives

Tableau 5d : synthèse des modèles

| Var | Nom | coefficient | MRO ordinaire | MRF sélection avant | MRB sélection arrière |
|-----|-----------------------|--------------------------------|------------------|---------------------------|-----------------------------|
| X0 | GENERAL intercepte | b0 | 41,4087 | 41,3312 | 36,8479 |
| X1 | CRIM | b1 | -0,1251 | -0,1252 | |
| X2 | NOX | b2 | -21,1251 | -20,8566 | -19,5103 |
| X3 | AGE | b3 | 0,0040 | | |
| X4 | DIS | b4 | -1,2916 | -1,3098 | -1,2817 |
| X5 | RM | b5 | 3,9002 | 3,9241 | 4,1822 |
| X6 | LSTAT | b6 | -0,5639 | -0,5589 | -0,6041 |
| X7 | RAD | b7 | 0,2685 | 0,2680 | |
| X8 | CHAS | b8 | 2,7103 | 2,7263 | |
| X9 | INDUS | b9 | -0,0022 | | |
| X10 | TAX | b10 | -0,0102 | -0,0102 | |
| X11 | PT | b11 | -0,9708 | -0,9673 | -0,9297 |
| | | SS resid résiduelle | 9826,13 | 9827,976 | 10580,48 |
| | | MSE = σ^2 (ANOVA) | 25,003 | 24,88095 | 26,517 |
| | | R ² | 0,71963493 | 0,719582 | 0,69811140 |
| | | R ² ajusté | ,71178756 | 0,713193 | 0,69432833 |

5e)

Comparez les prédictions des 3 modèles sur l'ensemble test T constitué des 101 observations. Choisir le meilleur modèle selon des critères; préciser la nature de ces critères.

| Var | Nom | coefficient | MRO ordinaire | MRF sélection avant | MRB sélection arrière |
|-----|-----------------------|--------------------------------|------------------|---------------------------|-----------------------------|
| X0 | GENERAL intercepte | b0 | 37,30622 | 33,70878 | 27,50599 |
| X1 | CRIM | b1 | 0,15954 | 0,12549 | |
| X2 | NOX | b2 | -4,99851 | | |
| X3 | AGE | b3 | -0,03647 | 0,028814 | |
| X4 | DIS | b4 | -0,92637 | 0,393779 | |
| X5 | RM | b5 | 4,16260 | 4,23967 | 4,35412 |
| X6 | LSTAT | b6 | -0,50708 | -0,49156 | -0,45228 |
| X7 | RAD | b7 | 0,05491 | | |
| X8 | CHAS | b8 | 3,83516 | 4,09416 | |
| X9 | INDUS | b9 | -0,04094 | 0,100942 | |
| X10 | TAX | b10 | -0,00483 | | |
| X11 | PT | b11 | -1,35093 | -1,32519 | -1,45293 |
| | | SS resid résiduelle | 1508,168 | 1522,431 | 1769,708 |
| | | MSE = σ^2 (ANOVA) | 16,94570 | 1654,82 | 18,244 |
| | | R ² | 0,802669 | 0,80080241 | 0,76844831 |
| | | R ² ajusté | 0,778279 | 0,78348088 | 0,76128691 |
| | | F | 32,910728 | 46,232 | 107,30 |

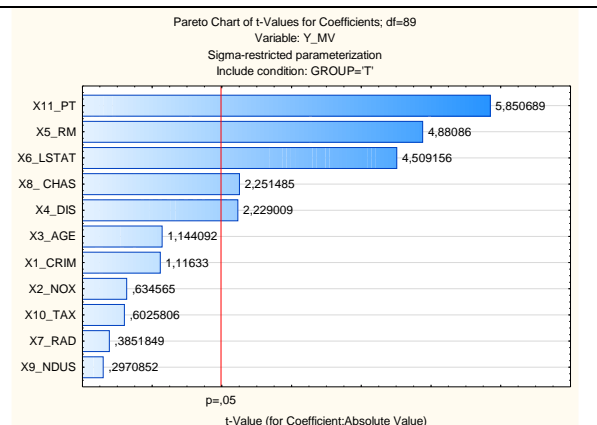
MOR

Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 R= .89591780 R²= .80266870 Adjusted R²= .77627944
 F(11,89)=32.911 p<0.0000 Std Error of estimate: 4.1165
 Include condition: GROUP=T

| N=101 | b ^a | Std Err. of b ^a | b | Std Err. of b | t(89) | p-value |
|-----------|----------------|----------------------------|-----------|---------------|----------|----------|
| Intercept | | | 37.30622 | 9.758018 | 3.82314 | 0.000244 |
| X1_CRIM | 0.105387 | 0.094405 | 0.105387 | 0.142915 | 1.11633 | 0.267286 |
| X2_NOX | -0.062573 | 0.098607 | -0.062573 | 0.142915 | -0.43456 | 0.527340 |
| X3_AGE | -0.118951 | 0.103970 | -0.118951 | 0.14409 | -0.83456 | 0.255653 |
| X4_DIS | -0.217179 | 0.097433 | -0.217179 | 0.145595 | -1.49409 | 0.028331 |
| X5_RM | 0.332211 | 0.068064 | 0.332211 | 0.085284 | 3.88086 | 0.000005 |
| X6_LSTAT | -0.414082 | 0.091831 | -0.414082 | 0.112456 | -3.68116 | 0.000020 |
| X7_RAD | 0.056561 | 0.146842 | 0.056561 | 0.142562 | 0.39718 | 0.701019 |
| X8_CHAS | -0.111972 | 0.049732 | -0.111972 | 0.073389 | -1.51449 | 0.026816 |
| X9_NDUS | -0.031720 | 0.106772 | -0.031720 | 0.137812 | -0.22970 | 0.767094 |
| X10_TAX | -0.094350 | 0.156577 | -0.094350 | 0.080811 | -1.16633 | 0.548320 |
| X11_PT | -0.346350 | 0.059198 | -0.346350 | 0.230901 | -1.49409 | 0.000000 |

Variables currently in the Equation: DV: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 Include condition: GROUP=T

| Variable | b ^a in | Partial Cor. | Semipart Cor. | Tolerance | R-square | t(89) | p-value |
|----------|-------------------|--------------|---------------|-----------|----------|----------|----------|
| X1_CRIM | 0.105387 | 0.117511 | 0.052565 | 0.248781 | 0.751219 | 1.11633 | 0.267286 |
| X2_NOX | -0.062573 | -0.067112 | -0.029880 | 0.228027 | 0.771973 | -0.43456 | 0.527340 |
| X3_AGE | -0.118951 | -0.120391 | -0.053872 | 0.205112 | 0.794888 | -0.83456 | 0.255653 |
| X4_DIS | -0.217179 | -0.229943 | -0.104958 | 0.233557 | 0.766443 | -1.49409 | 0.028331 |
| X5_RM | 0.332211 | 0.459513 | 0.229826 | 0.478597 | 0.521403 | 3.88086 | 0.000005 |
| X6_LSTAT | -0.414082 | -0.431242 | -0.212324 | 0.262920 | 0.737080 | -3.68116 | 0.000020 |
| X7_RAD | 0.056561 | 0.040796 | 0.018137 | 0.102826 | 0.897174 | 0.39718 | 0.701019 |
| X8_CHAS | -0.111972 | 0.232138 | 0.106016 | 0.896453 | 0.103547 | -1.51449 | 0.026816 |
| X9_NDUS | -0.031720 | -0.031475 | -0.013989 | 0.194487 | 0.805513 | -0.22970 | 0.767094 |
| X10_TAX | -0.094350 | -0.063744 | -0.028374 | 0.090437 | 0.909563 | -1.16633 | 0.548320 |
| X11_PT | -0.346350 | -0.527045 | -0.275493 | 0.632689 | 0.367311 | -1.49409 | 0.000000 |



Collinearity statistics for terms in the equation (BostonHousing.sta in 2018-MTH302-Devois-Data)
 Sigma-restricted parameterization
 Include condition: GROUP=T

| Effect | Tolerance | Variance Inflation | R square | Y_MV Beta in | Y_MV Partial | Y_MV Semi-par | Y_MV t | Y_MV p |
|----------|-----------|--------------------|-----------|--------------|--------------|---------------|-----------|-----------|
| X1_CRIM | 0.2487806 | 4.019606 | 0.7512194 | 0.1053870 | 0.1175109 | 0.0525649 | 1.116330 | 0.2672859 |
| X2_NOX | 0.2280270 | 4.385446 | 0.7719730 | -0.0625729 | -0.0671121 | -0.0298799 | -0.434565 | 0.5273400 |
| X3_AGE | 0.2051118 | 4.875389 | 0.7948882 | -0.1189511 | -0.1203914 | -0.0538721 | -0.834565 | 0.2556532 |
| X4_DIS | 0.2335566 | 4.281617 | 0.7664434 | -0.2171794 | -0.2299432 | -0.1049578 | -1.494092 | 0.0283314 |
| X5_RM | 0.4785972 | 2.089440 | 0.5214028 | 0.3322112 | 0.4595131 | 0.2298261 | 3.880860 | 0.0000046 |
| X6_LSTAT | 0.2629196 | 3.803444 | 0.7370804 | -0.4140824 | -0.4312416 | -0.2123236 | -3.681156 | 0.0000198 |
| X7_RAD | 0.1028263 | 9.725135 | 0.8971737 | 0.0565614 | 0.0407965 | 0.0181373 | 0.395185 | 0.7010194 |
| X8_CHAS | 0.8964528 | 1.115508 | 0.1035472 | 0.1119717 | 0.2321376 | 0.1060162 | 2.251485 | 0.0268164 |
| X9_NDUS | 0.1944867 | 5.141739 | 0.8055133 | -0.0317204 | -0.0314754 | -0.0139889 | -0.229708 | 0.7670937 |
| X10_TAX | 0.0904374 | 11.057378 | 0.9095626 | -0.0943505 | -0.0637435 | -0.0283738 | -1.166330 | 0.5483200 |
| X11_PT | 0.6326885 | 1.580557 | 0.3673115 | -0.3463499 | -0.5270451 | -0.2754926 | -1.494090 | 0.0000001 |

D'après le graphe et p-valeu on remarque qu'il y a des variables qui sont non significatifs, et on a aussi max VIF=11,05 donc on a le problème de multicollinéarité. Le modèle ne satisfaisant pas

MRE

Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 R= .89487564 R²= .80080241 Adjusted R²= .76348088
 F(8,92)=46.232 p<0.0000 Std Error of estimate: 4.0679
 Include condition: GROUP=T

| N=101 | b ^a | Std Err. of b ^a | b | Std Err. of b | t(92) | p-value |
|-----------|----------------|----------------------------|-----------|---------------|----------|----------|
| Intercept | | | 33.70878 | 7.932084 | 4.24968 | 0.000051 |
| X5_RM | 0.338362 | 0.065560 | 0.338362 | 0.0821465 | 4.11611 | 0.000001 |
| X11_PT | -0.339752 | 0.055712 | -0.339752 | 0.217304 | -1.56383 | 0.000000 |
| X6_LSTAT | -0.401407 | 0.089531 | -0.401407 | 0.109639 | -3.66432 | 0.000021 |
| X8_CHAS | 0.119534 | 0.047801 | 0.119534 | 0.072227 | 1.64067 | 0.014166 |
| X1_CRIM | 0.082898 | 0.066917 | 0.082898 | 0.101302 | 0.817382 | 0.218563 |
| X4_DIS | -0.223459 | 0.092319 | -0.223459 | 0.139779 | -1.59937 | 0.017463 |
| X3_AGE | -0.151513 | 0.093971 | -0.151513 | 0.104646 | -1.44714 | 0.150315 |
| X9_NDUS | -0.095175 | 0.078206 | -0.095175 | 0.100942 | -0.94367 | 0.226729 |

Analysis of Variance: DV: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 Include condition: GROUP=T

| Effect | Sums of Squares | df | Mean Squares | F | p-value |
|----------|-----------------|----|--------------|----------|----------|
| Regress. | 6120.385 | 8 | 765.0486 | 46.23162 | 0.000000 |
| Residual | 1522.431 | 92 | 16.5482 | | |
| Total | 7642.820 | | | | |

D'après les tableaux on remarque que la méthode de Régression avec Forward Stepwise a éliminé (négligé) quelques variables qui ne sont pas significatives mais il n'a pas éliminé toutes les variables puisque la méthode de Forward Stepwise ne vérifie pas l'états de toutes les variables après la sélection chaque step.

MRB

Regression Summary for Dependent Variable: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 R= .87661183 R²= .76844831 Adjusted R²= .76126691
 F(3,97)=107.30 p<0.0000 Std Error of estimate: 4.2713
 Include condition: GROUP=T

| N=101 | b ^a | Std Err. of b ^a | b | Std Err. of b | t(97) | p-value |
|-----------|----------------|----------------------------|-----------|---------------|----------|----------|
| Intercept | | | 27.50599 | 7.645390 | 3.59772 | 0.000507 |
| X5_RM | 0.347496 | 0.066119 | 0.347496 | 0.0828469 | 4.19562 | 0.000001 |
| X6_LSTAT | -0.369334 | 0.064799 | -0.369334 | 0.079352 | -4.65396 | 0.000000 |
| X11_PT | -0.372500 | 0.054104 | -0.372500 | 0.211030 | -1.76542 | 0.000000 |

Analysis of Variance: DV: Y_MV (BostonHousing.sta in Workbook1_(Recovered))
 Include condition: GROUP=T

| Effect | Sums of Squares | df | Mean Squares | F | p-value |
|----------|-----------------|----|--------------|----------|----------|
| Regress. | 5873.112 | 3 | 1957.704 | 107.3043 | 0.000000 |
| Residual | 1769.708 | 97 | 18.244 | | |
| Total | 7642.820 | | | | |

D'après les tableaux on remarque que la méthode de Régression avec Backward Stepwise a éliminé les variables qui ne sont pas significatives et résoudre le problème de multicollinéarité.

Backward Stepwise vérifie l'états de toutes les variables après la sélection pour chaque step

D'après les résultats précédentes on remarque que la méthode de Régression avec Backward Stepwise plus rapide et efficace et contient moins des étapes par rapport Régression avec Forward Stepwise et Régression Ordinaire.

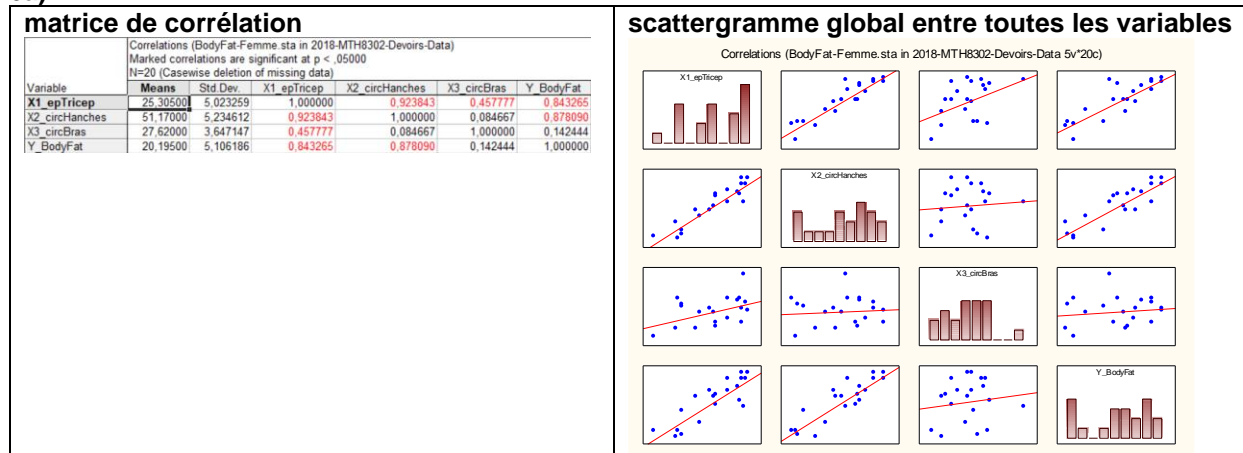
Et on remarque que avec la méthode Backward Stepwise on a éliminé le problème des variables non significatifs qui apparaissent dans le modèle de Régression Ordinaire et Régression avec Forward Stepwise, et aussi on remarque qu'il a éliminé le problème de multicollinéarité, Et aussi R^2_{adj} légèrement inférieur à R^2 .

No 6 Étude d'un modèle de régression multiple problématique

Données = BodyFat-F.sta

Réponse

6a)

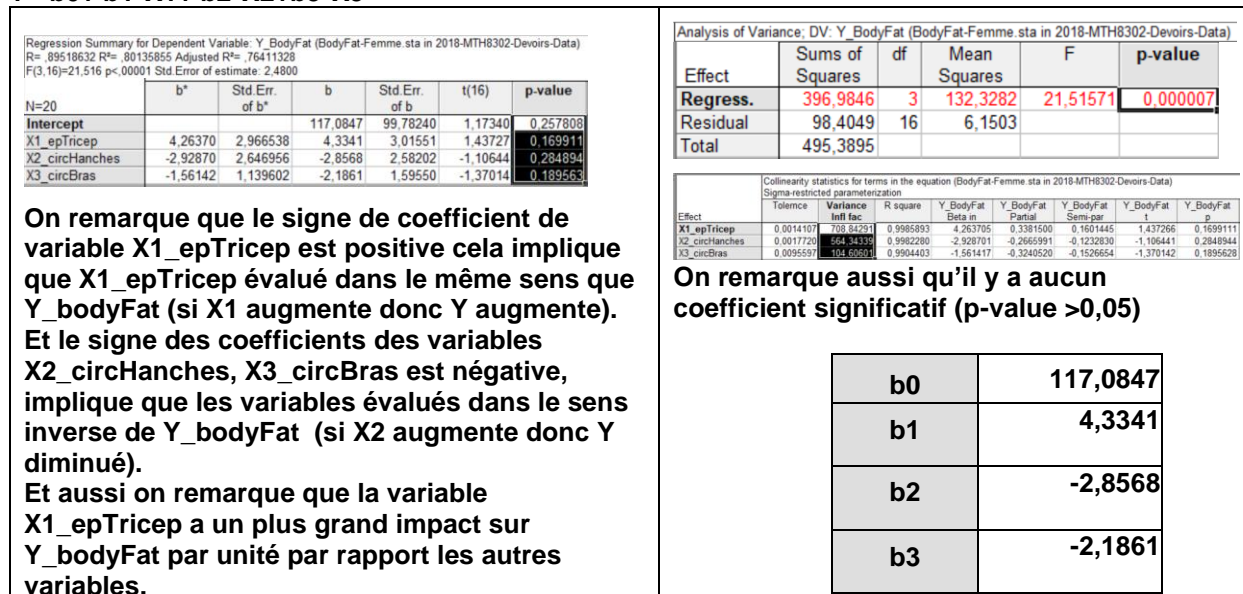


D'après la matrice de corrélation et le scattergramme on remarque que La corrélation est positive entre les variables et est valide on remarque qu'il y a une forte corrélation entre (X2 et X1 ,...) ce qui implique qu'il y a un problème de multicollinéarité

6b)

Modèle de régression multiple ordinaire (MRO) entre Y et X1, X2, X3.

$$Y = b_0 + b_1 \cdot X_1 + b_2 \cdot X_2 + b_3 \cdot X_3$$



Le modèle ne satisfaisant pas à cause de présentent un problème de multicollinéarité entre les variable Xi (max VIF=708,84 >10)

Pour obtenir un modèle plus satisfaisant on utilise des méthodes pour remédier aux problèmes de multicollinéarité comme la régression en composantes principales (ACP), la régression PLS (« Partial Least Square »), la régression RIDGE, la méthode de sélection de variables...

Et pour cet exercice on utilise les 2 méthodes suivant pour résoudre le problème de multicollinéarité : Mod1 : régression RIDGE et Mod2 : régression en composantes principales

6c)

Mod1 : régression RIDGE

K=0

Ridge Regression Summary for Dependent Variable: Y_BodyFat (BodyFat-Femme.sta in 2018-MTH8302-Devoirs-Data)
t=0.0000 R= .89518632 R²= .80135855 Adjusted R²= .76411328
F(3,16)=21.516 p<.00001 Std Error of estimate: 2.4800

| | b* | Std. Err. of b* | b | Std. Err. of b | t(16) | p-value |
|----------------|----------|-----------------|----------|----------------|----------|----------|
| Intercept | | | 117.0847 | 99.78240 | 1.17340 | 0.257808 |
| X1_epTricep | 4.26370 | 2.966538 | 4.3341 | 3.01551 | 1.43727 | 0.169911 |
| X2_circHanches | -2.92870 | 2.646956 | -2.8568 | 2.58202 | -1.10644 | 0.284894 |
| X3_circBras | -1.56142 | 1.139602 | -2.1861 | 1.59550 | -1.37014 | 0.189563 |

K=0,1

Ridge Regression Summary for Dependent Variable: Y_BodyFat (BodyFat-Femme.sta in 2018-MTH8302-Devoirs-Data)
t=0.0000 R= .89518632 R²= .80135855 Adjusted R²= .76411328
F(3,16)=21.516 p<.00001 Std Error of estimate: 2.4800

| | b* | Std. Err. of b* | b | Std. Err. of b | t(16) | p-value |
|----------------|----------|-----------------|----------|----------------|----------|----------|
| Intercept | | | 117.0847 | 99.78240 | 1.17340 | 0.257808 |
| X1_epTricep | 4.26370 | 2.966538 | 4.3341 | 3.01551 | 1.43727 | 0.169911 |
| X2_circHanches | -2.92870 | 2.646956 | -2.8568 | 2.58202 | -1.10644 | 0.284894 |
| X3_circBras | -1.56142 | 1.139602 | -2.1861 | 1.59550 | -1.37014 | 0.189563 |

...

La méthode RIDGE [cours ch04-Multiple2 p21]

Régression « RIDGE »

Méthode Hoerli et Kennard

$b_{MC} = \beta = (X'X)^{-1} X'y$ estimateur de moindres carrés (MC)

nouvelle classe d'estimateurs appelés **RIDGE (R)**

$b_R(k) = (X'X + kI)^{-1} X'y$ $k > 0$ $b_R(0) = b_{MC}$

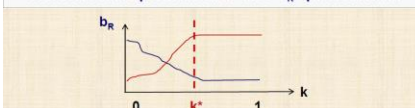
k : paramètre biaisant - à déterminer

$b_R = c b_{MC}$ $0 < c < 1$ **interprétation**

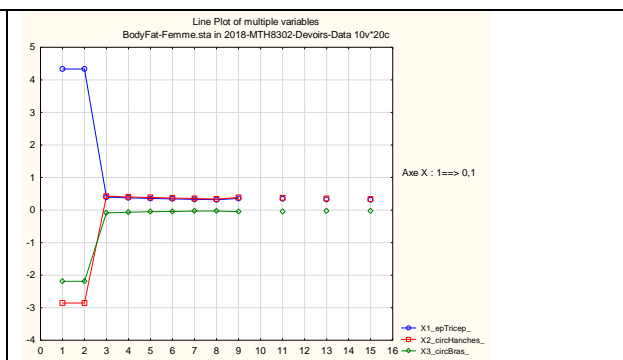
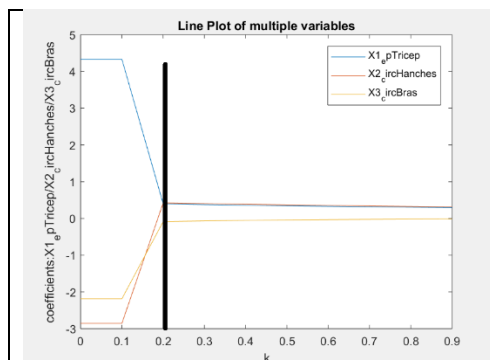
- estimateur biaisé rétréci d'erreur minimale
- évite l'instabilité des coefficients b_{MC} en multicollinéarité

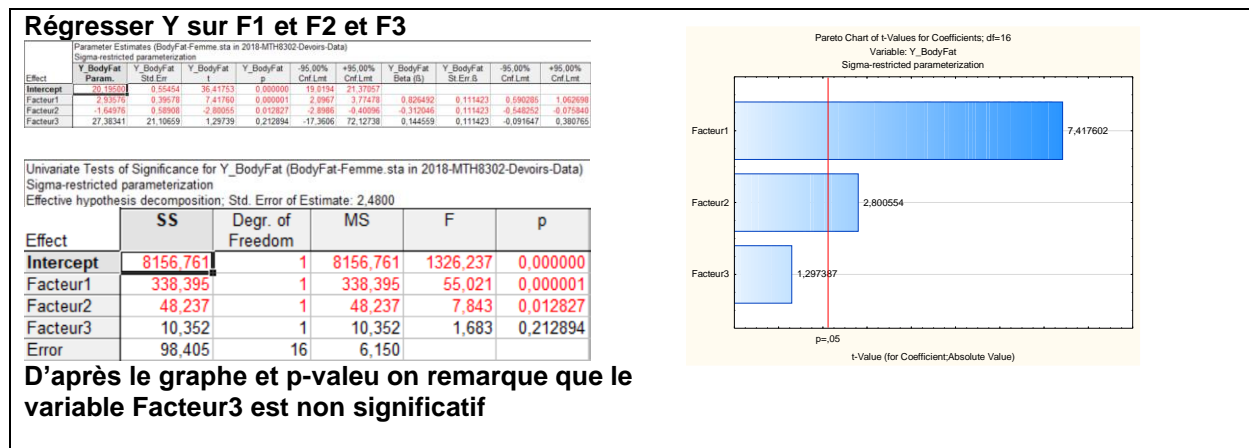
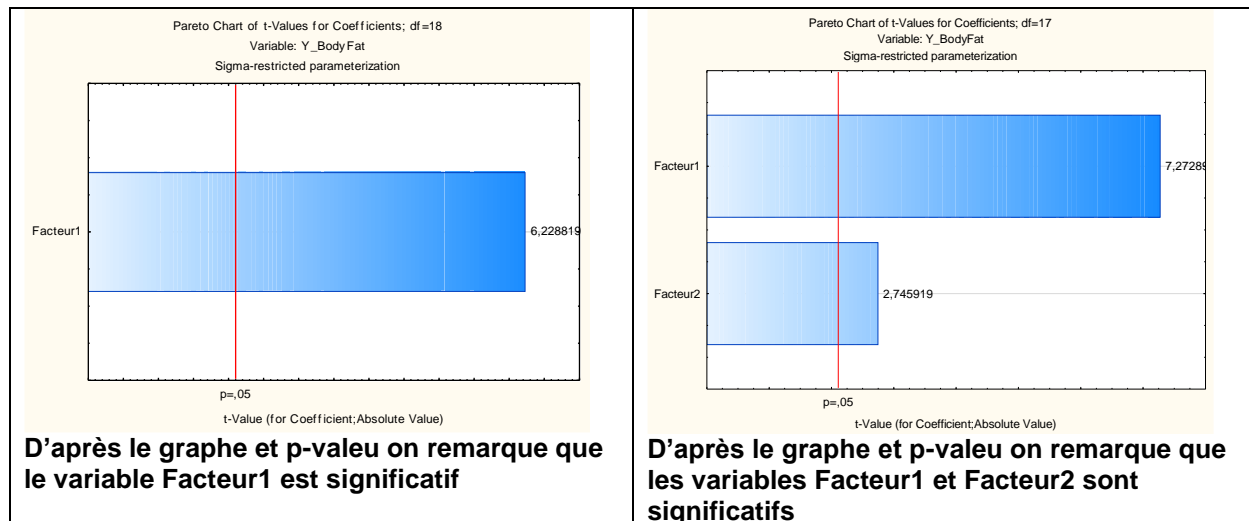
Choix de k graphique de b_R en fonction de k («Ridge Trace»)

recherche d'une petite valeur k^* avec b_R quasi constants



| k | X1 | X2 | X3 |
|-----|---------|---------|----------|
| 0,0 | 4,3341 | -2,8568 | -2,1861 |
| 0,1 | 4,3341 | -2,8568 | -2,1861 |
| 0,2 | 0,39789 | 0,42405 | -0,08581 |
| 0,3 | 0,37644 | 0,40521 | -0,06704 |
| 0,4 | 0,35877 | 0,38685 | -0,05270 |
| 0,5 | 0,34330 | 0,36976 | -0,04130 |
| 0,6 | 0,32939 | 0,35401 | -0,03205 |
| 0,7 | 0,31673 | 0,33953 | -0,02443 |





Et après on fait la Comparaisant entre Yprédit et Yobservé et on modifie l'étape de régression et prend le meilleur choix.

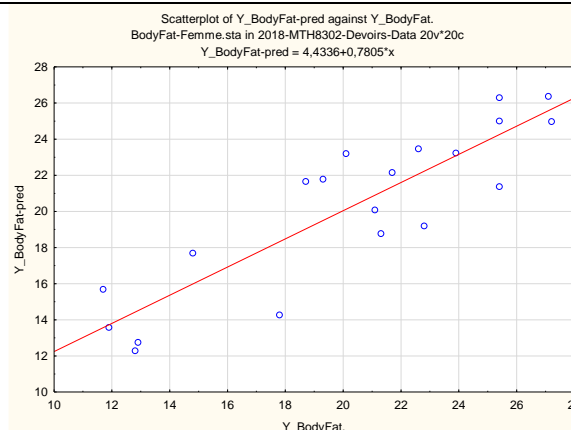
Et a la fin nous avons trouvé que $Y = b_0 + b_1 \cdot F_1 + b_2 \cdot F_2$ c'est le meilleur choix

Pour régresser Y sur F1 et F2

Nous avons trouvé

$$Y = b_0 + b_1 \cdot F_1 + b_2 \cdot F_2$$

| Var | Coefficient | MR |
|-----|--------------------------------|--------------------|
| F0 | b0 | 20,19500 |
| F1 | b1 | 2,93576 |
| F2 | b2 | -1,64976 |
| | SS resid résiduelle | 108,7572 28,12% |
| | MSE = σ^2 (ANOVA) | 6,397480 |
| | R ² | 0,780461 |
| | R ² ajusté | 0,754633 |



D'après le graphe et le tableau on remarque que on a une bonne corrélation entre les facteurs et Y_pred (R²ajusté=0,75). Et aussi qu'il y a 28% de variabilité de variable qui n'est pas expliqué par notre modèle

| Variable | Factor 1 | Factor 2 | Factor 3 |
|----------------|----------|-----------|-----------|
| X1_epTricep | 0,998641 | -0,048393 | 0,019342 |
| X2_circHanches | 0,904817 | -0,425451 | -0,017255 |
| X3_circBras | 0,500494 | 0,865708 | -0,007399 |
| *Y_BodyFat | 0,826492 | -0,312046 | 0,144559 |

$$Y_{\text{BodyFat_pred}} = b_0 + b_1 \cdot F_1 + b_2 \cdot F_2 = b_0 + (b_1 \cdot 0,998 - b_2 \cdot 0,048) \cdot X1_{\text{epTricep_cr}} + (b_1 \cdot 0,904 - b_2 \cdot 0,425) \cdot X2_{\text{cirHanches}} + (b_1 \cdot 0,5 - b_2 \cdot 0,312) \cdot X3_{\text{circBras}}$$

6d)

| | Mod1 | Mod2 |
|-----------------------|----------|----------|
| R ² | 0,686220 | 0,780461 |
| R ² ajusté | 0,627386 | 0,754633 |
| SSresid(%) | 45,72 | 28,12 |

D'après la valeur de R²ajusté et SSresid(%) on remarque que le Mod2 (régression en composantes principales) le meilleur choix. Et on remarque aussi qu'il y a des variables non significatives dans le Mod1 (régression RIDGE) par contre le Mod2 où toutes ces variables sont significatives.

No 7 Étude de prédiction d'activité biologique : modélisation PLS

Données = Penta.sta

Réponse

7a)

M1 : Modèle PLS (modèle avec toutes les composantes)

Predictor weights (Penta.sta in 2018-MTH8302-Devoirs-Data)

Responses: Y_logRAI

Options: NO-INTERCEPT AUTOSCALE

Include condition: CLASSE='entraînement'

| | "S1" | "L1" | "P1" | "S2" | "L2" | "P2" | "S3" | "L3" | "P3" | "S4" |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| Comp 1 | -0.151641 | 0.065681 | -0.169113 | 0.121527 | 0.071134 | 0.061188 | -0.424809 | 0.653701 | 0.285046 | -0.233037 |
| Comp 2 | -0.153810 | 0.021349 | -0.404829 | 0.229796 | 0.071941 | 0.250337 | -0.108375 | 0.325934 | 0.425376 | 0.414880 |
| Comp 3 | 0.130046 | -0.113962 | 0.001702 | -0.125617 | -0.246332 | 0.283583 | -0.340124 | 0.286179 | -0.092797 | 0.081971 |
| Comp 4 | 0.118332 | 0.025651 | 0.263189 | -0.108027 | -0.094759 | 0.006084 | -0.288687 | 0.372122 | -0.179677 | 0.161053 |
| Comp 5 | -0.045688 | 0.228215 | 0.310581 | 0.080716 | 0.204099 | -0.399491 | 0.007050 | 0.538996 | -0.199468 | 0.159731 |
| Comp 6 | 0.061752 | 0.080174 | 0.285662 | -0.288948 | -0.115836 | -0.263954 | 0.331086 | 0.719173 | -0.212239 | 0.191910 |
| Comp 7 | 0.110030 | 0.026349 | 0.087751 | 0.049820 | 0.036832 | 0.013176 | 0.359529 | 0.696413 | -0.556903 | 0.145516 |
| Comp 8 | -0.108092 | 0.122918 | -0.069392 | 0.000680 | -0.042247 | 0.086976 | 0.501340 | 0.605843 | -0.529243 | 0.035130 |
| Comp 9 | 0.109628 | -0.145432 | -0.178732 | 0.057439 | 0.059079 | -0.020814 | 0.615097 | 0.559399 | -0.282889 | -0.062998 |
| Comp 10 | 0.319291 | 0.378214 | 0.489696 | 0.231033 | -0.035344 | -0.037074 | 0.322718 | 0.073336 | 0.361522 | -0.064807 |
| Comp 11 | 0.679530 | 0.665916 | -0.307054 | 0.004732 | 0.001234 | 0.005058 | -0.003764 | -0.017029 | 0.005456 | -0.006089 |

| | "S5" | "L5" | "P5" |
|--|-----------|-----------|-----------|
| | 0.056905 | 0.056905 | -0.056905 |
| | -0.008193 | -0.008193 | 0.008193 |
| | -0.259076 | -0.259076 | 0.259076 |
| | 0.156438 | 0.156438 | -0.156438 |
| | -0.115688 | -0.115688 | 0.115688 |
| | -0.048718 | -0.048718 | 0.048718 |
| | -0.014484 | -0.014484 | 0.014484 |
| | 0.009015 | 0.009015 | -0.009015 |
| | 0.021920 | 0.021920 | -0.021920 |
| | -0.017691 | -0.017691 | 0.017691 |
| | 0.001854 | 0.001854 | -0.001854 |

PLS regression coefficients (Penta.sta in 2018-MTH8302-Devoirs-Data)

Responses: Y_logRAI

Options: NO-INTERCEPT AUTOSCALE

Include condition: CLASSE='entraînement'

| | Intercept | "S1" | "L1" | "P1" | "S2" | "L2" | "P2" | "S3" | "L3" | "P3" | "S4" |
|----------|-----------|----------|----------|-----------|----------|----------|----------|-----------|----------|----------|-----------|
| Y_logRAI | 0.737328 | 0.126345 | 0.268079 | -0.442360 | 0.101932 | 0.058532 | 0.000840 | -0.049595 | 0.362006 | 0.117779 | -0.031100 |

| | "L4" | "P4" | "S5" | "L5" | "P5" |
|--|-----------|-----------|----------|----------|-----------|
| | -0.057821 | -0.328174 | 0.011786 | 0.090855 | -0.008913 |

Summary of PLS (Penta.sta in 2018-MTH8302-Devoirs-Data)

Responses: Y_logRAI

Options: NO-INTERCEPT AUTOSCALE

Include condition: CLASSE='entraînement'

| | Increase R² of Y | Average R² of Y | Increase R² of X | Average R² of X |
|---------|------------------|-----------------|------------------|-----------------|
| Comp 1 | 0.896399 | 0.896399 | 0.169014 | 0.169014 |
| Comp 2 | 0.078368 | 0.974767 | 0.127721 | 0.296735 |
| Comp 3 | 0.004636 | 0.979403 | 0.146554 | 0.443289 |
| Comp 4 | 0.002485 | 0.981889 | 0.118421 | 0.561710 |
| Comp 5 | 0.001494 | 0.983383 | 0.105894 | 0.667605 |
| Comp 6 | 0.002617 | 0.986001 | 0.051876 | 0.719481 |
| Comp 7 | 0.002428 | 0.988428 | 0.061873 | 0.781354 |
| Comp 8 | 0.001926 | 0.990354 | 0.072252 | 0.853606 |
| Comp 9 | 0.000725 | 0.991080 | 0.067285 | 0.920891 |
| Comp 10 | 0.000000 | 0.991080 | 0.079076 | 0.999967 |
| Comp 11 | 0.000099 | 0.991179 | 0.000033 | 1.000000 |

D'après les R2 de de Y(cumul) on remarque que nous avons une bonne corrélation entre les composantes et Y

Y_logRAI-pred=b0+b1*S1+...b15*P5

| Y_logRAI-OBSV | Y_logRAI-pred-com11 | écart |
|---------------|---------------------|----------|
| 0,00 | 0,10001 | -0,10001 |
| 0,28 | 0,32348 | -0,04348 |
| 0,20 | 0,09999 | 0,10001 |
| 0,51 | 0,32345 | 0,18655 |
| 0,11 | 0,11000 | 0,00000 |
| 2,73 | 2,60227 | 0,12773 |
| 0,18 | 0,18000 | 0,00000 |
| 1,53 | 1,58613 | -0,05613 |
| -0,10 | -0,10000 | 0,00000 |
| -0,52 | -0,52000 | 0,00000 |
| 0,40 | 0,40000 | 0,00000 |
| 0,30 | 0,30000 | 0,00000 |
| -1,00 | -1,00000 | 0,00000 |
| 1,57 | 1,71053 | -0,14053 |
| 0,59 | 0,66413 | -0,07413 |

Scatterplot of Y_logRAI_pred-com11 against Y_logRAI

Penta.sta in 2018-MTH8302-Devoirs-Data 21v*46c

Include condition: CLASSE='entraînement'

Y_logRAI_comp11 = 0,004+0.9912*x

On remarque que le modèle prédit a bien suivi notre modèle réel

7b)

modèle à 2 composantes semble un bon choix R2 de Y = 0,9747

Summary of PLS (Penta.sta in 2018-MTH8302-Devoirs-Data)
 Responses: Y_logRAI
 Options: NO-INTERCEPT AUTOSCALE
 Include condition: CLASSE='entraînement'

| | Increase R² of Y | Average R² of Y | Increase R² of X | Average R² of X |
|---------|---------------------|--------------------|---------------------|--------------------|
| Comp 1 | 0,896399 | 0,896399 | 0,169014 | 0,169014 |
| Comp 2 | 0,078368 | 0,974767 | 0,127721 | 0,296735 |
| Comp 3 | 0,004636 | 0,979403 | 0,146554 | 0,443289 |
| Comp 4 | 0,002485 | 0,981889 | 0,118421 | 0,561710 |
| Comp 5 | 0,001494 | 0,983383 | 0,105894 | 0,667605 |
| Comp 6 | 0,002617 | 0,986001 | 0,051876 | 0,719481 |
| Comp 7 | 0,002428 | 0,988428 | 0,061873 | 0,781354 |
| Comp 8 | 0,001926 | 0,990354 | 0,072252 | 0,853606 |
| Comp 9 | 0,000725 | 0,991080 | 0,067285 | 0,920891 |
| Comp 10 | 0,000000 | 0,991080 | 0,079076 | 0,999967 |
| Comp 11 | 0,000099 | 0,991179 | 0,000033 | 1,000000 |

On a basé sur le choix des composantes sur les deux critères suivants :

Nombre de composante petit et R2 de Y (cumul) max
 Et comme vous voyez dans le tableau nous remarquons que R2 de Y(cumul) de la composante2 augmente de 0,07 par rapport R2 de Y(cumul) de composante1, et après la composante2 on remarque que les autres composantes augmentent de façon presque constante.
 avec R2 de Y (cumul)=0,97 qui montre que notre modèle presque prédit tous les variables

M2 : Modèle PLS (modèle avec les 2 premières composantes)

Predictor weights (Penta.sta in 2018-MTH8302-Devoirs-Data)

Responses: Y_logRAI

Options: NO-INTERCEPT AUTOSCALE

Include condition: CLASSE='entraînement'

| | "S1" | "L1" | "P1" | "S2" | "L2" | "P2" | "S3" | "L3" | "P3" |
|---------|-----------|----------|-----------|----------|----------|----------|-----------|----------|----------|
| Compo 1 | -0,157641 | 0,085681 | -0,169313 | 0,121527 | 0,071134 | 0,065188 | -0,424809 | 0,653701 | 0,285046 |
| Compo 2 | -0,193810 | 0,021349 | -0,404829 | 0,229796 | 0,071941 | 0,250337 | -0,108375 | 0,325834 | 0,425376 |

| | "S4" | "L4" | "P4" | "S5" | "L5" | "P5" |
|-----------|-----------|-----------|-----------|-----------|-----------|------|
| -0,293407 | 0,298287 | -0,203095 | 0,056905 | 0,056905 | -0,056905 | |
| 0,414880 | -0,283113 | 0,355025 | -0,008193 | -0,008193 | 0,008193 | |

PLS regression coefficients (Penta.sta in 2018-MTH8302-Devoirs-Data)

Responses: Y_logRAI

Options: NO-INTERCEPT AUTOSCALE

Include condition: CLASSE='entraînement'

| | Intercept | "S1" | "L1" | "P1" | "S2" | "L2" | "P2" | "S3" | "L3" |
|----------|-----------|-----------|----------|-----------|----------|----------|----------|-----------|----------|
| Y_logRAI | -0,437289 | -0,083797 | 0,039170 | -0,296360 | 0,145064 | 0,085710 | 0,106487 | -0,106782 | 0,209140 |

| | "P3" | "S4" | "L4" | "P4" | "S5" | "L5" | "P5" |
|----------|-----------|----------|-----------|----------|----------|-----------|------|
| 0,515502 | -0,101102 | 0,117977 | -0,254619 | 0,028722 | 0,221406 | -0,021719 | |

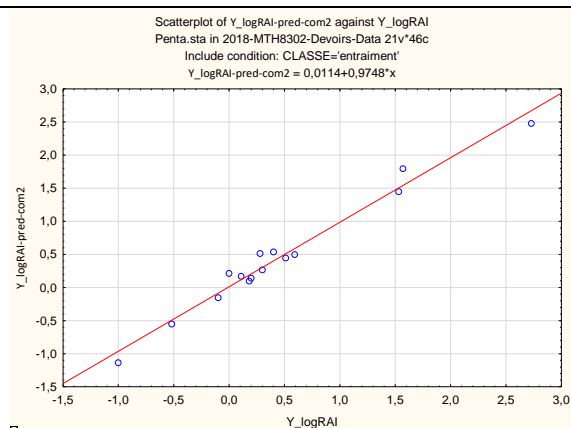
Y_logRAI-pred=b0+b1*S1+...b15*P5

Summary of PLS (Penta.sta in 2018-MTH8302-Devoirs-Data)
 Responses: Y_logRAI
 Options: NO-INTERCEPT AUTOSCALE
 Include condition: CLASSE='entraînement'

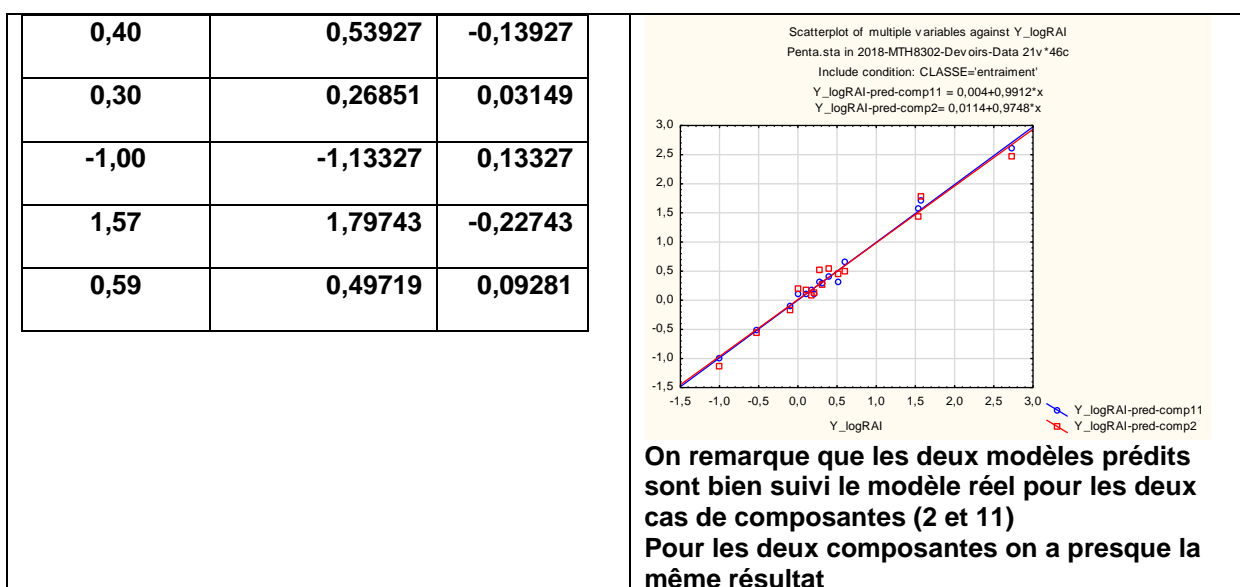
| | Increase R² of Y | Average R² of Y | Increase R² of X | Average R² of X |
|--------|---------------------|--------------------|---------------------|--------------------|
| Comp 1 | 0,896399 | 0,896399 | 0,169014 | 0,169014 |
| Comp 2 | 0,078368 | 0,974767 | 0,127721 | 0,296735 |

D'après les R2 de de Y(cumul) on remarque que nous avons une bonne corrélation entre les composantes et Y

| Y_logRAI-OBSV | Y_logRAI-pred-com2 | écart |
|---------------|--------------------|----------|
| 0,00 | 0,21319 | -0,21319 |
| 0,28 | 0,51533 | -0,23533 |
| 0,20 | 0,14380 | 0,05620 |
| 0,51 | 0,44595 | 0,06405 |
| 0,11 | 0,17156 | -0,06156 |
| 2,73 | 2,48083 | 0,24917 |
| 0,18 | 0,09637 | 0,08363 |
| 1,53 | 1,44762 | 0,08238 |
| -0,10 | -0,15456 | 0,05456 |
| -0,52 | -0,54923 | 0,02923 |



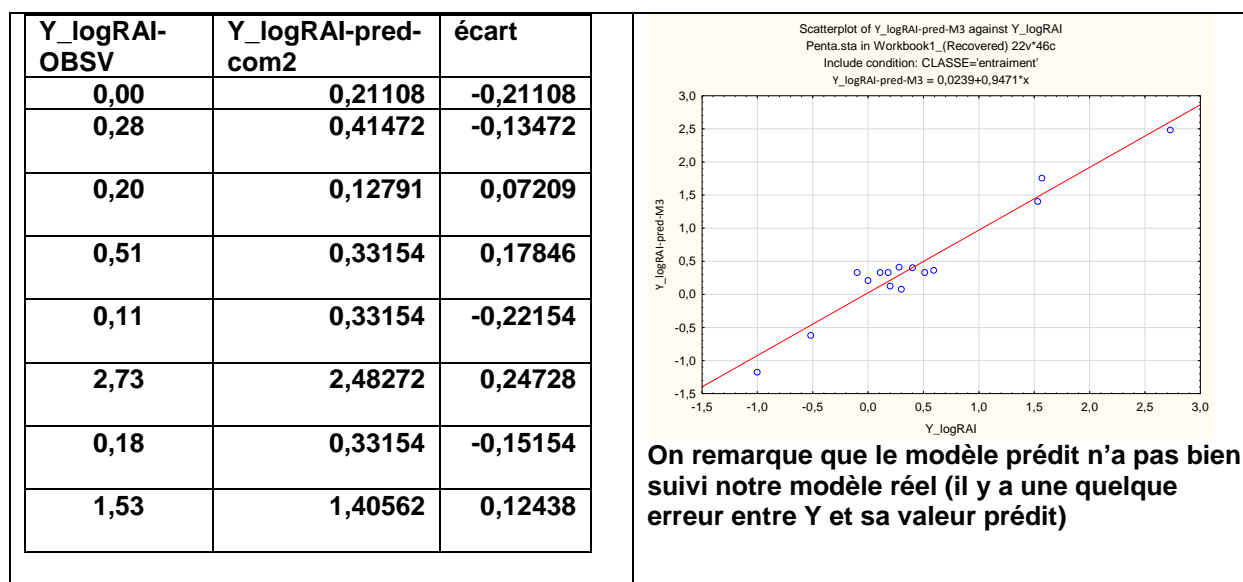
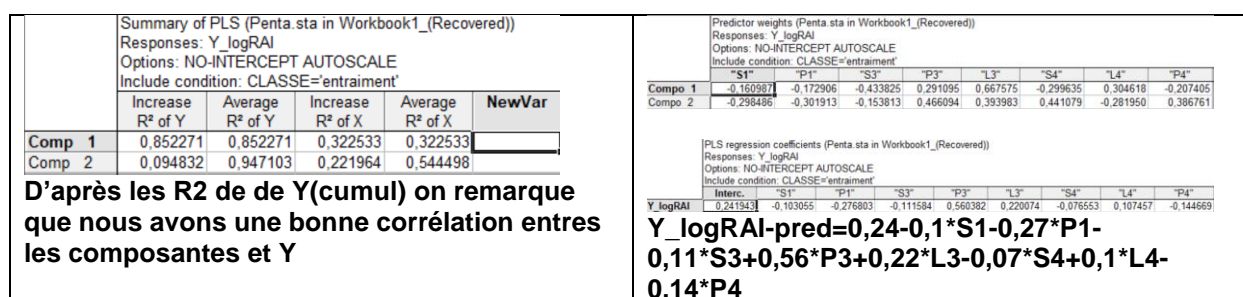
On remarque que le modèle prédit a bien suivi notre modèle réel (il y a une linéarisation entre Y et sa valeur prédit)



7c)

Avec les variables S1 P1 S3 P3 L3 S4 L4 P4 on a trouvé R2 de Y pour composant2 égale 0,947 est ca signifie que avec c'est variable on a exprimé 95% de valeur prédit de Y et a cause de ca on peut éliminé les autres variables qui ne contribué pas avec les valeurs prédites

M3 : Modèle PLS (modèle avec les 2 premières composantes et les variables S1 P1 S3 P3 L3 S4 L4 P4)

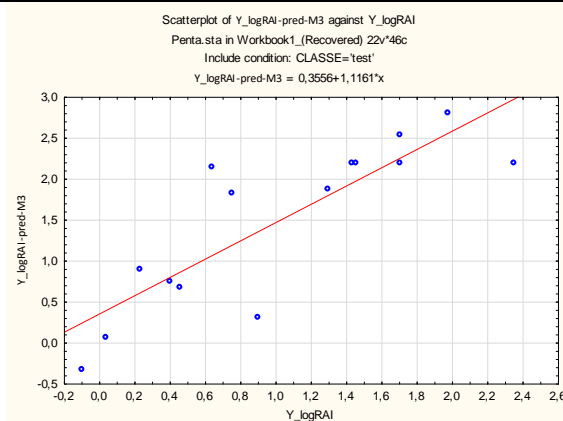


| | | |
|-------|----------|----------|
| -0,10 | 0,33154 | -0,43154 |
| -0,52 | -0,61608 | 0,09608 |
| 0,40 | 0,40105 | -0,00105 |
| 0,30 | 0,08113 | 0,21887 |
| -1,00 | -1,17345 | 0,17345 |
| 1,57 | 1,75619 | -0,18619 |
| 0,59 | 0,36295 | 0,22705 |

7d)

$$Y_{\log RAI-pred} = 0,24 - 0,1 * S1 - 0,27 * P1 - 0,11 * S3 + 0,56 * P3 + 0,22 * L3 - 0,07 * S4 + 0,1 * L4 - 0,14 * P4$$

| Y_logRAI-OBSV | Y_logRAI-pred-com2 | écart |
|---------------|--------------------|-------|
| -0,10 | -0,327235642 | 0,23 |
| 0,46 | 0,686589345 | -0,23 |
| 0,75 | 1,82393778 | -1,07 |
| 1,43 | 2,20 | -0,78 |
| 1,45 | 2,20 | -0,76 |
| 1,71 | 2,20 | -0,50 |
| 0,04 | 0,05 | -0,01 |
| 0,23 | 0,89 | -0,66 |
| 1,30 | 1,86 | -0,57 |
| 2,35 | 2,20 | 0,14 |
| 1,98 | 2,80 | -0,82 |
| 1,71 | 2,53 | -0,82 |
| 0,90 | 0,30 | 0,60 |
| 0,64 | 2,14 | -1,51 |
| 0,40 | 0,75 | -0,35 |
| | -1,89 | 1,89 |



On remarque que le modèle prédit n'a pas bien suivi notre modèle réel (il y a une grand erreur entre Y et sa valeur prédit)

Pour que la régression soit plus rapide et efficace on élimine les variables qui ne contribué pas avec les valeurs prédites et on prend le nombre de composante petit et R2 de Y (cumul) max, avec cela nous obtenons le même résultat

Conclusion générale

Dans cet devoir nous avons traité des données réelles ou on a utilisé la régression multiple ordinaire et Régression avec Forward Stepwise et Régression avec Backward Stepwise. Et on a traité le problème de multicollinéarité ou on a utilisé les méthodes (RIDGE, ACP) pour résoudre le problème. Et on a bien étudié la régression PLS