
TRAVAIL PRATIQUE 2

TRAVAIL PRÉSENTÉ À
M. THIERRY DUCHESNE

DANS LE CADRE DU COURS
THÉORIE ET APPLICATIONS DES MÉTHODES DE RÉGRESSION
STT-7125

RÉALISÉ PAR L'ÉQUIPE 21 :
ALEXANDRE LEPAGE
& AMEDEO ZITO

LE 17 DÉCEMBRE 2020



UNIVERSITÉ
LAVAL

FACULTÉ DES SCIENCES ET DE GÉNIE
ÉCOLE D'ACTUARIAT
UNIVERSITÉ LAVAL

Table des matières

1	Modèle linéaire mixte pour les résultats en mathématique	1
1a)	Exclusion de la variable <code>meanses</code>	1
1b)	Inclusion de la variable <code>meanses</code>	4
2	Modèle linéaire mixte pour la grandeur de jeunes filles	5
3	GEE pour le nombre d'auto-administrations de doses analgésiques	8
3a)	Modèle linéaire généralisé	8
3b)	Prédiction pour la population	8
A	Graphiques	9
A.1	Question 1	9
A.2	Question 2	15
A.3	Question 3	16

Liste des illustrations

1	Résidus en fonction des valeurs prédites pour le modèle (1).	9
2	Résidus en fonction des valeurs prédites pour le modèle (3).	9
3	Résidus en fonction des valeurs prédites pour le modèle (10).	10
4	Splines réalisés sur les variables <code>homework</code> et <code>ratio</code> lors de l'entraînement d'un GAM.	10
5	Visualisation des Splines suite à l'entraînement d'un GAM utilisant (2).	10
6	Résidus studentisés en fonction de l'index des observations pour le modèle 3.	11
7	Résidus studentisés en fonction de l'index des observations pour le modèle 10.	11
8	Résidus studentisés en fonction des différentes variables explicatives du modèle 3.	12
9	Résidus studentisés en fonction des différentes variables explicatives du modèle 10.	13
10	Sortie R de la fonction <code>summary</code> pour le modèle (7).	14
11	Sortie R de la fonction <code>summary</code> pour le modèle (11).	14
12	Relation de la grandeur en fonction de l'âge pour chacune des jeunes filles.	15
13	Graphiques de résidus générés à partir du modèle (15).	15
14	Sortie R de la fonction <code>summary</code> pour le modèle (17).	16

Introduction

Les modèles linéaires (LM) et modèles linéaires généralisés (GLM) sont des outils fort utiles pour modéliser toute sorte de phénomènes et sont largement utilisés dans le milieu statistique. Cependant, ces modèles s'appuient sur l'hypothèse que les observations Y_1, \dots, Y_n , $n > 0$ servant à les entraîner sont indépendantes ; laquelle n'est pas toujours réaliste selon le contexte.

Les modèles linéaires mixtes (LMM) permettent donc d'insérer une structure de dépendance entre les observations d'un LM. Du côté des GLM, il est possible d'effectuer un ajustement au modèle pour que la matrice de variance de celui-ci puisse tenir compte d'éventuelles covariances. De tels modèles ajustés sont appelés les GEE (*Generalized Estimating Equation*).

L'objet de ce travail pratique est de mettre en pratique ces deux modèles. Ainsi, les questions 1 et 2 abordent le sujet des modèles linéaires mixtes tandis que la question 3 aborde le sujet des GEE.

1 Modèle linéaire mixte pour les résultats en mathématique

Pour la première question de ce travail, on regarde les données d'un sous-ensemble des étudiants de 8ème année ayant participé au *National Educationnal Longitudinal Study* de 1988. L'objectif de cette étude est de voir comment les résultats en mathématiques varient en fonction du nombre d'heures de travail à la maison (variable `homework` dans la base de données). Dans ce cas-ci, on a que la variable endogène Y_{ij} correspond au résultat de l'examen de mathématique de l'étudiant j , appartenant à l'école i , $i = 1, \dots, 23$, $j = 1, \dots, n_i$, où n_i correspond au nombre d'étudiants appartenant à la i -ème grappe. Dans la base de données, cette variable est désignée comme `math`.

Comme la variable explicative `meanses`, correspondant au statut socio-économique moyen des étudiants de l'école, est fortement corrélée avec l'école d'origine des étudiants (variable de grappe), alors on est intéressé de voir quelle différence il y aurait entre un LMM avec et sans cette variable et de voir si le besoin d'effets aléatoires dans le modèle persiste si on ajoute cette dernière.

1a) Exclusion de la variable meanses

Pour débiter l'entraînement d'un modèle, la première étape est de considérer un LM, d'évaluer ses résidus pour voir si les postulats sont respectés et de prendre action autrement. À noter que les LMM permettent de régler les problèmes d'auto-corrélation des résidus et, dans certains cas, de régler l'hétéroscédasticité.

Entraînement d'un modèle linéaire

On considère le modèle linéaire suivant :

$$Y_{ij} = \beta_0 + \beta_1 \text{homework} + \beta_2 \text{white} + \beta_3 \text{ratio} + \epsilon_{ij}. \quad (1)$$

Une analyse de multicollinéarité réalisée conformément à la méthodologie décrite dans le travail pratique 1 du présent cours ne soulève aucun problème. En revanche, lorsque l'on regarde l'illustration 1, on voit que les résidus ont une légère tendance descendante. Le postulat de linéarité n'est donc pas respecté. Pour remédier à ce problème, on regarde les splines générés par un modèle additif généralisé (GAM) à l'aide de la fonction `gam` du *package* du même nom. Ceux-ci sont présentés dans l'illustration 4.

On voit donc dans l'illustration 4a qu'il est possible de passer une droite dans l'intervalle de confiance entourant le spline. Pour cette variable, aucune transformation n'est donc nécessaire. Pour ce qui est de la variable `ratio`, il est impossible de passer une telle droite. Conséquemment, il faudrait ajouter un terme de deuxième degré sur la variable `ratio` qui aurait été centrée et réduite au préalable. Pour se faire, on pose

$$\text{ratio}^* = \frac{\text{ratio} - 18}{\sigma_{\text{ratio}}} \quad \text{et} \quad \text{ratio2} = (\text{ratio}^*)^2.$$

Le modèle linéaire devient alors

$$Y_{ij} = \beta_0 + \beta_1 \text{homework} + \beta_2 \text{white} + \beta_3 \text{ratio}^* + \beta_4 \text{ratio2} + \epsilon_{ij}. \quad (2)$$

Si on refait l'exercice du GAM, on trouve l'illustration 5. Dans celle-ci, on remarque qu'aucune transformation additionnelle n'est nécessaire. Si on se fie à la statistique F produite par la fonction `summary` en R, alors on trouve que la seule non-linéarité qui soit significative au seuil de 5% est celle de la variable `homework`. Cependant, par soucis de simplicité, celle-ci ne sera pas modifiée. Par la suite, on peut tenter d'ajouter des interactions. Au seuil de 5%, les interactions ajoutées sont `white:ratio` et `white:ratio2`, ce qui permet d'obtenir le modèle (3)

$$Y_{ij} = \beta_0 + \beta_1 \text{homework} + \beta_2 \text{white} + \beta_3 \text{ratio}^* + \beta_4 \text{ratio2} + \beta_5 (\text{white} : \text{ratio}^*) + \beta_6 (\text{white} : \text{ratio2}) + \epsilon_{ij}. \quad (3)$$

L'illustration 2 permet de voir que la linéarité semble meilleure, mais que le problème est désormais au niveau de l'hétéroscédasticité des résidus. Or, un LMM peut aider à traiter ce genre de problème. Du point de vue de l'auto-corrélation des observations du jeu de données, l'illustration 6 permet de voir que, selon l'école d'appartenance des élèves, les résidus du LM ne sont pas identiquement distribués. En effet, selon la grappe, on voit que les résidus ont une variance et une moyenne qui peut différer. Cette observation vient donc légitimer l'utilisation d'un LMM.

Entraînement d'un modèle linéaire mixte

Maintenant qu'un modèle linéaire a été entraîné et que l'on a observé la nécessité d'y inclure des effets aléatoires pour tenir compte de la corrélation qui existe entre les élèves d'une même école, il est temps d'entraîner un LMM.

Pour se faire, la première étape est de faire les graphiques des résidus en fonction des différentes variables pour voir sur lesquelles d'entre elles il serait intéressant appliquer un effet aléatoire. Bien que cela soit difficile à voir, l'illustration 8 montre que la variable la plus susceptible d'avoir un effet aléatoire est `homework` puisqu'elle est celle ayant le plus de volatilité dans la distribution des résidus selon les valeurs qu'elle peut prendre. Ainsi, avec la fonction `lmer` du package R `lme4`, on entraîne le modèle (4) avec les structures de variance VC¹ pour la variance des résidus et UN², de même que UN(1)³ pour la variance des effets aléatoires.

$$Y_{ij} = \beta_0 + \gamma_{i0} + (\beta_1 + \gamma_{i1}) \text{homework} + \beta_2 \text{white} + \beta_3 \text{ratio}^* + \beta_4 \text{ratio2} + \beta_5 (\text{white} : \text{ratio}^*) + \beta_6 (\text{white} : \text{ratio2}) + \epsilon_{ij}. \quad (4)$$

À noter que l'ajout de trop d'effets aléatoires dans le modèle testé peut entraîner de l'instabilité numérique lors de l'entraînement de celui-ci. C'est pourquoi on se limite à deux effets aléatoires dans le modèle (4).

En ce qui attrait à la structure de variance des résidus de type CS⁴, comme la fonction `lmer` ne permet pas de l'utiliser, on peut faire appel à la fonction `lme` du package `nlme`. Pour ce qui est de la structure AR(1)⁵, celle-ci n'a que peu de sens dans ce contexte puisque les observations (les élèves d'une même école) ne peuvent pas être ordonnées selon un ordre chronologique ou spatial. Pour cette raison, on ne considérera pas cette dernière.

Comme les trois modèles testés possèdent tous la même composante fixe ($\mathbf{X}\beta$), alors on peut comparer les log-vraisemblances de même que les AIC. Le tableau 1 présente donc l'AIC pour chacun des modèles testés.

1. *Variance Components* : indépendance entre les termes de résidus.
2. *Unstructured* : il existe une corrélation entre les effets aléatoires.
3. *Diagonales principales* : les effets aléatoires sont indépendants l'un de l'autre.
4. *Compound symmetry* : la corrélation entre les résidus est la même partout.
5. *Auto-régression d'ordre 1* : la corrélation diminue selon un aspect d'éloignement (généralement pour les observations qui sont étudiées à travers le temps ou l'espace).

Var(ϵ)	Var(γ)	dl	AIC
VC	UN	11.00	3630.72
VC	UN(1)	10.00	3658.96
CS	UN	12.00	3621.86

Tableau 1 – AIC des trois modèles testés en fonction de (4) avec le nombre de degrés de liberté dl associé à chacun d’eux.

Avec le tableau 1, on voit que la structure de variance qui minimise l’AIC est CS/UN. Cependant, avec la fonction `summary` de R, on voit que le coefficient de corrélation des résidus d’une même classe est de 5.126496×10^{-18} , ce qui est très près de zéro. On peut donc simplifier le modèle et prendre la structure VC/UN. Puis on voit que la corrélation entre les effets aléatoires γ_{i0} et γ_{i1} est de -0.89, ce qui confirme qu’il existe un lien de dépendance significatif entre ces variables aléatoires et que la structure de variance UN est approprié. En somme, la corrélation entre les étudiants d’une même école est négligeable et il existe un lien de dépendance significatif entre les effets aléatoires du modèle.

Après avoir sélectionné les structures de variance du LMM, il faut tester si les effets aléatoires du modèle (4) sont nécessaires. Pour se faire, il s’agit de procéder à un test du ratio des vraisemblances. Soit les hypothèses de test suivantes :

H_0 : Le modèle simple est suffisant ;

H_1 : Le modèle complet représente mieux les données.

Soit l_0 et l_1 , la log-vraisemblance sous H_0 et celle sous H_1 . On définit Δ_{dl} comme la différence du nombre de paramètres entre les deux modèles. Le calcul de la p -value du test est effectué avec (5).

$$p\text{-value} = 0.5 [2 - \mathbb{P}(\chi_{\Delta_{dl}-1}^2 > \xi) - \mathbb{P}(\chi_{\Delta_{dl}}^2 > \xi)] , \quad (5)$$

On applique ainsi (5) pour évaluer si l’effet aléatoire γ_{i1} est significatif et on trouve une statistique de test de 92.92 avec $\Delta_{dl} = 2$, ce qui permet de calculer un seuil observé de 0. Conséquemment, on rejette fortement H_0 et on conserve l’effet aléatoire γ_{i1} . De plus, puisque γ_{i1} est conservé, on ne peut retirer l’ordonnée à l’origine aléatoire. Le modèle obtenu suite à cette étape de construction du LMM correspond ainsi à (6).

$$Y_{ij} = \beta_0 + \gamma_{i0} + (\beta_1 + \gamma_{i1})\text{homework} + \beta_2\text{white} + \beta_3\text{ratio}^* + \beta_4\text{ratio2} + \beta_5(\text{white} : \text{ratio}^*) + \beta_6(\text{white} : \text{ratio2}) + \epsilon_{ij}. \quad (6)$$

Il ne reste plus qu’à sélectionner les effets fixes. Pour se faire, on utilise le test de Wald de type III utilisé par la fonction `Anova` du *package* `car`. On remarque alors que la variable `ratio*` possède un seuil de test de 15.69%. Cependant, comme on ne peut la retirer sans avoir retiré les variables dépendantes d’elle au préalable, c.-à-d. `white:ratio2`, `white:ratio*` et `ratio2`, on ne peut pas l’enlever. Conséquemment, on va commencer par retirer l’interaction `white:ratio2` avant de réeffectuer le test. Puis, on retire aussi `white:ratio*` puisque la variable `ratio*` n’est toujours pas significative au seuil de 5%. On fait de même avec `ratio2` pour finalement retirer `ratio*`. On trouve ainsi le modèle final (7).

$$Y_{ij} = \beta_0 + \gamma_{i0} + (\beta_1 + \gamma_{i1})\text{homework} + \beta_2\text{white} + \epsilon_{ij}. \quad (7)$$

Avec la fonction `summary` de R, on obtient les résultats présentés dans l’illustration 10. D’une part, on a les effets fixes pour lesquels un intervalle de confiance à 95% est calculé dans le tableau 2.

	Estimateur	Écart-type	IC 95%	
β_0	44.02	1.83	40.42	47.62
β_1	1.90	0.92	0.11	3.70
β_2	3.30	0.98	1.38	5.22

Tableau 2 – Estimateurs des poids pour les effets fixes du LMM ainsi que leurs intervalles de confiance à 95%.

D'autre part, on a

$$D_i = \text{Var}(\gamma_i) = \begin{bmatrix} 58.20797 & -27.01225 \\ -27.01225 & 17.25707 \end{bmatrix}, i = 1, \dots, 23 \text{ et } \mathbf{D} = \begin{bmatrix} D_1 & 0 & \dots & 0 \\ 0 & D_2 & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & D_{23} \end{bmatrix} \quad (8)$$

De plus,

$$\mathbf{V} = \text{Var}(\epsilon) = 52.66 \mathbf{I}_{n \times n}, \quad n = \sum_{i=1}^{23} n_i = 519. \quad (9)$$

Discussion

Avec le tableau 2, on voit que, toute autre chose étant égale, chaque heure de travail supplémentaire à la maison contribue à augmenter la note moyenne d'un étudiant pour son examen de mathématique de 1.90%. Cet effet est significatif puisque l'intervalle de confiance à 95% de l'estimateur n'inclut pas la valeur 0. Par ailleurs, comme on a pu l'observer lors de l'étape du test des effets aléatoires, l'effet du nombre d'heures de travail à la maison peut varier d'une école à l'autre.

1b) Inclusion de la variable meanses

Comme pour la partie 1a), on commence par entraîner un modèle linéaire mixte et on procède de façon similaire pour trouver le modèle (10).

$$Y_{ij} = \beta_0 + \beta_1 \text{meanses} + \beta_2 \text{homework} + \beta_3 \text{white} + \beta_4 \text{ratio}^* + \beta_5 \text{ratio2} \\ + \beta_6 (\text{white} : \text{ratio}) + \beta_7 (\text{white} : \text{ratio2}) + \beta_8 (\text{meanses} : \text{white}) + \beta_9 (\text{meanses} : \text{ratio2}) + \epsilon_{ij}. \quad (10)$$

À noter que l'interaction des variables **homework** et **meanses** n'est pas significatif au seuil de 5%.

En comparant les illustrations 6 et 7, on voit que l'ajout de la variable **meanses** au modèle semble, a priori, régler le problème de corrélation entre les observations. En effet, on voit dans l'illustration 7 que les résidus de chacune des écoles semblent centrées à zéro. Cependant les variance varient encore quelque peu. Voyons maintenant si l'ajout d'effets aléatoires serait significatif.

Entraînement d'un modèle linéaire mixte

Pour débiter, l'illustration 9 montre que, encore une fois, seule la variable **homework** est susceptible de recevoir un effet aléatoire. Afin de confirmer cette observation, on peut entraîner un LMM ne comportant que deux effets aléatoires, soit une ordonnée à l'origine et un effet sur l'une des variables explicatives parmi **homework**, **meanses** et **ratio***. On teste ainsi chacune des variables avec les trois structures mentionnées ci-haut, soit VC/UN, VC/UN(1) et CS/UN. Il en découle que le modèle qui minimise l'AIC, est celui incluant un effet aléatoire à la variable **homework**. Par la suite, si on tente d'ajouter un troisième effet aléatoire, on obtient que les fonctions **lmer** et **lme** deviennent instables numériquement. On s'en tiendra donc à 2 effets. En ce qui attrait aux structures de variances, celle qui minimise l'AIC est la structure CS/UN. Cependant, comme dans la section 1a), on a un coefficient de corrélation pour la variance de ϵ_i qui est de 5.126496×10^{-18} . La dépendance entre les résidus d'une même grappe est donc négligeable et on peut simplifier le modèle en adoptant la structure VC/UN. Plus encore, la corrélation entre l'ordonnée à l'origine aléatoire et l'effet appliqué à la variable **homework** est de -0.91 confirmant ainsi que la structure UN est adéquate pour la variance de γ .

Par la suite, on effectue le test du ratio des vraisemblances dont le calcul du seuil observé est présenté en (5). On trouve ainsi une statistique de test de 90.95 avec $\Delta_{dl} = 2$, pour un seuil observé de 0. L'évidence est donc forte contre l'hypothèse nulle et on peut en conclure que l'effet aléatoire ajouté à la variable **homework** est significatif.

Pour ce qui est de la sélection des effets fixes, on a que la variable ayant le plus grand seuil observé avec le test de Wald de type III est **ratio***. Cependant, comme pour la section 1a), on doit gérer les variables d'ordre supérieur qui dépendent de celle-ci avant de pouvoir la retirer. On commence donc par retirer l'interaction **white:ratio2**. Puis, on refait le test pour retirer **ratio2**; ainsi de suite jusqu'à trouver le modèle (11) où tous les effets fixes sont significatifs au seuil de 5%.

$$Y_{ij} = \beta_0 + \gamma_{i0} + \beta_1 \text{meanses} + (\beta_2 + \gamma_{i2}) \text{homework} + \beta_3 \text{white} + \epsilon_{ij}. \quad (11)$$

Si on essaie d'intégrer l'interaction **meanses:homework**, on trouve un seuil de test de 0.7368; on ne l'inclue donc pas dans le modèle. La sortie R de la fonction **summary** appliquée sur le modèle ainsi obtenu est présentée dans l'illustration 11. Les effets fixes sont décrits dans le tableau 3 et les matrices de variances sont présentées dans (12) et (13), lesquelles sont exactement les mêmes que dans la section 1a).

	Estimateurs	Écart-types	IC 95%	
β_0	44.70	1.79	41.20	48.21
β_1	4.89	1.34	2.26	7.52
β_2	1.93	0.90	0.17	3.68
β_3	3.11	0.96	1.24	4.99

Tableau 3 – Estimateurs des poids pour les effets fixes du LMM (11) ainsi que leurs intervalles de confiance à 95%.

$$D_i = \text{Var}(\gamma_i) = \begin{bmatrix} 58.20797 & -27.01225 \\ -27.01225 & 17.25707 \end{bmatrix}, i = 1, \dots, 23 \text{ et } D = \begin{bmatrix} D_1 & 0 & \dots & 0 \\ 0 & D_2 & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & D_{23} \end{bmatrix} \quad (12)$$

$$V = \text{Var}(\epsilon) = 52.66 \mathbf{I}_{n \times n}, \quad n = \sum_{i=1}^{23} n_i = 519. \quad (13)$$

Discussion

Comme on a pu l'observer avec l'illustration 7, l'ajout de la variable **meanses** qui est très fortement corrélée avec les identifiants des écoles (les grappes) a réduit considérablement le besoin d'ajouter des effets aléatoires au modèle puisque les résidus sont maintenant centrés autour de zéro. Néanmoins, avec le test du ratio des vraisemblances, on a pu voir que l'effet aléatoire appliqué à la variable **homework**, de même que l'ordonnée à l'origine aléatoire, sont utiles.

Au final, on trouve que l'ajout d'une heure supplémentaire d'étude augmente, en moyenne, l'espérance de la note en mathématique de 1.925% et cet effet varie d'une école à l'autre (effet aléatoire).

2 Modèle linéaire mixte pour la grandeur de jeunes filles

Pour cette deuxième question, le jeu de données à l'étude présente 20 courbes de la croissance de jeunes filles mesurées annuellement entre les âges 6 à 10 ans. Celui-ci a été publié par Gildstein (1979). Dans ce cas-ci, la variable endogène Y_{ij} correspond à la taille de la i -me fille lors de sa j -ème mesure à l'âge $5 + j$, $i = 1, \dots, 20$, $j = 1, \dots, 5$.

Entraînement d'un modèle linéaire

Afin de voir si la relation qui existe entre l'âge des petites filles et leur grandeur est linéaire, on regarde l'illustration 12. Comme celle-ci l'est effectivement, on entraîne le modèle (14).

$$Y_{ij} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{group}_2 + \beta_3 \text{group}_3 + \beta_4 (\text{group}_2 : \text{age}) + \beta_5 (\text{group}_3 : \text{age}) + \epsilon_{ij}. \quad (14)$$

Comme (14) possède un facteur d'inflation de la variance généralisé ($\text{GVIF}_j^{1/(2p_j)}$) supérieur à $\sqrt{10} = 3.16$, on est en présence de multicolinéarité. Pour remédier à ce problème, on peut simplement tronquer la variable **age** de la manière suivante :

$$\text{temps} = \text{age} - 6.$$

Le modèle (14) devient alors (15).

$$Y_{ij} = \beta_0 + \beta_1 \text{temps} + \beta_2 \text{group}_2 + \beta_3 \text{group}_3 + \beta_4 (\text{group}_2 : \text{temps}) + \beta_5 (\text{group}_3 : \text{temps}) + \epsilon_{ij}. \quad (15)$$

Avec ce dernier, on calcule les résidus studentisés de manière à générer l'illustrations 13. Dans un premier temps, on remarque avec l'illustration 13a que les résidus ne sont pas tous centrés autour de zéro. Dépendamment de la fillette, ceux-ci ont une moyenne qui diffère grandement, ce qui laisse présager une corrélation entre les observations d'une même fillette. Cela suggère qu'un LMM pourrait régler le problème d'auto-corrélation des résidus.

Entraînement d'un modèle linéaire mixte

Dans un deuxième temps, on remarque avec les illustrations 13a et 13b que les deux graphiques sont pratiquement identiques, laissant présager qu'un effet aléatoire sur la variable **temps** n'aurait aucune incidence sur les résidus. Plus encore, avec l'illustration 13c, on voit que les résidus varient énormément selon la valeur de la variable **group**. Cela pourrait expliquer en partie les ordonnées à l'origine des résidus qui diffèrent dans 13a.

Voyons maintenant si ces observations s'avèrent réalistes en entraînant un LMM avec les structures de variances VC/UN, CS/UN et AR(1)/UN. À noter que la structure AR(1) pour la variance des résidus est particulièrement intéressante dans ce contexte puisque les observations d'une même fillette peuvent être ordonnées chronologiquement. À cet effet, avec ce dernier, l'ajout de la variable **temps** dans les effets aléatoires engendre des problèmes de convergence avec la fonction **lme**. Conséquemment, le modèle entraîné à cette étape consiste en (16).

$$Y_{ij} = \beta_0 + \gamma_{i0} + \beta_1 \text{temps} + (\beta_2 + \gamma_{i2}) \text{group}_2 + (\beta_3 + \gamma_{i3}) \text{group}_3 + \beta_4 (\text{group}_2 : \text{temps}) + \beta_5 (\text{group}_3 : \text{temps}) + \epsilon_{ij}. \quad (16)$$

L'AIC calculé pour chacun des modèles entraînés est présenté dans le tableau 4

Var(ϵ)	Var(γ)	dl	AIC
VC	UN	13.00	356.72
CS	UN	14.00	360.87
AR(1)	UN	14.00	336.09

Tableau 4 – AIC des trois modèles testés en fonction de (16) avec le nombre de degrés de liberté *dl* associé à chacun d'eux.

On voit avec le tableau 4 que la structure de variance la plus appropriée selon le critère de l'AIC pour le modèle (16) est AR(1)/UN. Plus encore, le coefficient de corrélation liant les résidus d'une même fillette est de 0.9041, ce qui est hautement significatif. De plus, la matrice des coefficients de corrélation des effets aléatoires s'exprime comme

$$\rho(\gamma_i) = \begin{bmatrix} 1 & 0.380 & -0.147 \\ 0.380 & 1 & -0.813 \\ -0.147 & -0.813 & 1 \end{bmatrix}, \quad i = 1, \dots, 20.$$

Ainsi, la structure de variance non structurée (UN) est justifiée pour les effets aléatoires puisque les coefficients de corrélations sont significativement différents de zéro.

Si on fait un test du ratio des vraisemblances pour l'effet aléatoire appliqué sur la variable **group**, on trouve une statistique de test de 1.39 pour $\Delta_{dl} = 5$, ce qui donne un seuil observé de 0.885. Ainsi, on ne peut rejeter l'hypothèse nulle que l'effet aléatoire associé à la variable **group** n'est pas significatif et on

peut le retirer. Si on refait le test sur l'ordonnée à l'origine aléatoire, on trouve une statistique de 185.54 pour $\Delta_{dl} = 2$, ce qui donne un seuil observé de 0. On ne peut donc pas retirer l'ordonnée à l'origine aléatoire. Mentionnons que pour ce dernier test, comme on compare un modèle linéaire (H_0) à un modèle mixte (H_1), on ne peut utiliser la méthode REML pour calculer la log-vraisemblance de ce dernier lors du test puisque le modèle linéaire associé à l'hypothèse nulle utilise le maximum de vraisemblance.

Au niveau des effets fixes, le test de Wald de type III indique que tous les effets fixes sont significatifs au seuil de 5%. On obtient donc le modèle final (17)

$$Y_{ij} = \beta_0 + \gamma_{i0} + \beta_1 \text{temps} + \beta_2 \text{group}_2 + \beta_3 \text{group}_3 + \beta_4 (\text{group}_2 : \text{temps}) + \beta_5 (\text{group}_3 : \text{temps}) + \epsilon_{ij}. \quad (17)$$

La sortie R de la fonction `summary` appliquée sur le modèle ainsi obtenu est présentée dans l'illustration 14. Les effets fixes sont décrits dans le tableau 5 et les matrices de variances sont présentées dans (18) et (19).

	Estimateurs	Écarts-types	IC 95%	
β_0	112.57	1.22	110.18	114.96
β_1	3.70	1.66	0.44	6.96
β_2	7.79	1.66	4.54	11.05
β_3	5.29	0.19	4.92	5.65
β_4	0.26	0.25	-0.23	0.76
β_5	0.87	0.25	0.37	1.37

Tableau 5 – Estimateurs des poids pour les effets fixes du LMM (17) ainsi que leurs intervalles de confiance à 95%.

$$\mathbf{D} = 5.088543 \times 10^{-6} \mathbf{I}_{20 \times 20} \quad (18)$$

et

$$\mathbf{V}_i = \begin{bmatrix} 2.988 & 2.838 & 2.696 & 2.560 & 2.432 \\ 2.838 & 2.988 & 2.838 & 2.696 & 2.560 \\ 2.696 & 2.838 & 2.988 & 2.838 & 2.696 \\ 2.560 & 2.696 & 2.838 & 2.988 & 2.838 \\ 2.432 & 2.560 & 2.696 & 2.838 & 2.988 \end{bmatrix}, \quad i = 1, \dots, 20, \quad \mathbf{V} = \begin{bmatrix} V_1 & 0 & \dots & 0 \\ 0 & V_2 & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & V_{20} \end{bmatrix} \quad (19)$$

Discussion

L'interprétation des résultats obtenus dans le tableau 5 est synthétisée dans le tableau 6.

Taille de la mère	Taille de la fille à 6 ans	Taux de croissance annuel	Grandeur de la fille à 10 ans
petite	112.57	3.70	127.37
moyenne	120.36	3.96	136.20
grande	117.86	4.57	136.14

Tableau 6 – Mesures moyennes(en cm) pour une petite fille qui est née en 1973 selon la grandeur de la mère.

Avec la matrice \mathbf{D} obtenue en (18), on voit que l'effet aléatoire tend à être dégénérée puisque sa variance tend vers zéro. En somme, on peut résumer l'effet aléatoire γ_{i0} comme une constante, soit 0. Néanmoins, le modèle mixte est utile puisqu'il fait un ajustement sur la variance des prévisions pour tenir compte de la corrélation entre les mesures d'une même petite fille. On a donc $\mathbf{Y}_i \sim N(\mathbf{X}'_i \boldsymbol{\beta}, \mathbf{V}_i)$ où les composantes de la matrice $\boldsymbol{\beta}$ sont définis dans le tableau 5 et \mathbf{V}_i est défini en (19).

À la lumière de ces résultats, on peut répondre à la question de Goldstein (1979) en affirmant que oui, la croissance des filles est liée à la taille de la mère. L'interaction `group:temps` qui est significative au seuil de 0.001731 en atteste et les résultats du tableau 6 l'illustre bien.

3 GEE pour le nombre d'auto-administrations de doses analgésiques

Cette question a comme objectif d'utiliser un GEE afin de modéliser une étude sur effet d'addiction de doses analgésique. Les données contiennent $i = 1, \dots, 65$ patients sur $j = 1, \dots, 12$ période de 4 heures. Les patients sont ensuite divisés en deux groupes, une avec une dose de 1 mg, soit $x_i = 0$ et une avec une dose de 2 mg de morphine, soit $x_i = 1$. On définit donc un modèle de dénombrement complet comme

$$Y_{ij} = \beta_0 + \beta_1 x_i + \beta_2 t_j + \beta_3 x_i t_j,$$

avec Y_{ij} le nombre de dose auto-administré dans l'intervalle de temps t_j .

3a) Modèle linéaire généralisé

Dans ce contexte on peut utiliser un modèle de régression Poisson. Sachant que les Y_{ij} pour tout j sont dépendants entre eux, on ne peut pas simplement utiliser un GLM. Afin de tenir compte de cette dépendance à l'intérieur des grappes, on va utiliser un modèle linéaire généralisé.

On a testé plusieurs différentes matrices de corrélation de travail : La non structurée (UN), la auto-régressive (AR1), la structure échangeable et finalement la structure d'indépendance. On suppose que la drogue a un impact à court terme. Les administrations de drogue les plus récentes sont les plus importantes et sont donc plus corrélées. De ce fait, l'impact se dégrade dans le temps. Selon cette hypothèse le modèle avec AR(1) serait le plus approprié. On obtient pour la première ligne de la matrice

$$R_i(\alpha) = \begin{pmatrix} 1.000000000 & 0.516425829 & 0.266695637 & 0.137728516 & 0.071126563 & 0.03673159 & 0.01896914 \\ 0.009796156 & 0.005058988 & 0.002612592 & 0.001349210 & 0.0006967669 & & \end{pmatrix}.$$

Ce qui implique un $\alpha = 0.516425829$.

Il y a que deux variables exogènes dans le modèle, on fait des tests chi-carrés pour valider leur importance statistique. On commence avec l'interaction x_i et t_j . La valeur-p avec 1 degré de liberté est 0.13. Donc on ne peut pas rejeter H_0 et donc on écarte l'interaction. On fait pareille pour β_1 et β_2 . On obtient 0.015 et 0.0 respectivement, ce qui conclut qu'on garde les deux variables, puisqu'on rejette H_0 . Le modèle final est donc

$$Y_{ij} = \beta_0 + \beta_1 x_i + \beta_2 t_j.$$

Les paramètres sont

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} 2.05378263 \\ 0.30742284 \\ -0.06091784 \end{pmatrix}$$

Il semble qu'augmenter la dose à 2mg a une relation positive avec le nombre de dose administré. Le temps a une relation inverse et diminue l'espérance de 5% à chaque 4 heures écoulés (j augmente).

3b) Prédiction pour la population

On fait une estimation ponctuelle pour $x_i = 1$ et $t_j = 5$. On obtient

$$\eta = c\beta = (1 \quad 1 \quad 5) \begin{pmatrix} 2.05378263 \\ 0.30742284 \\ -0.06091784 \end{pmatrix} = 2.056616,$$

avec un intervalle de confiance 95% de

$$(1.87615, 2.237077).$$

Pour $E[Y_{i5}]$ on arrive à 7.819466 dans un intervalle de confiance de (6.528357, 9.365917). L'effet de la dose pour la population est de 0.3074228 avec un intervalle de confiance de (0.06038578, 0.55445990). Ainsi augmenter la dose à 2mg augmente en moyenne le nombre d'administrations de 35% avec IC entre (6%, 74%).

A Graphiques

A.1 Question 1

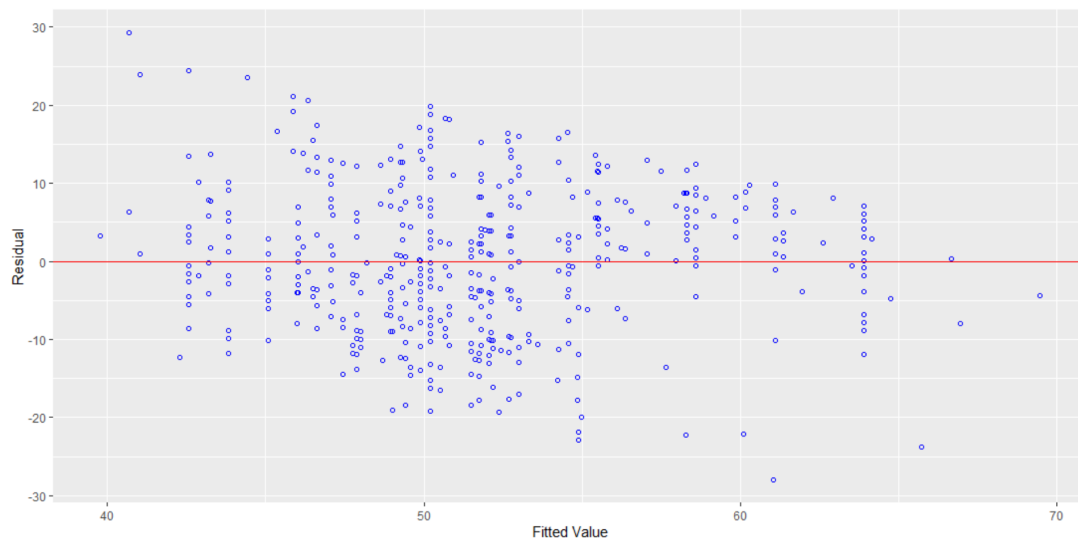


Illustration 1 – Résidus en fonction des valeurs prédites pour le modèle (1).

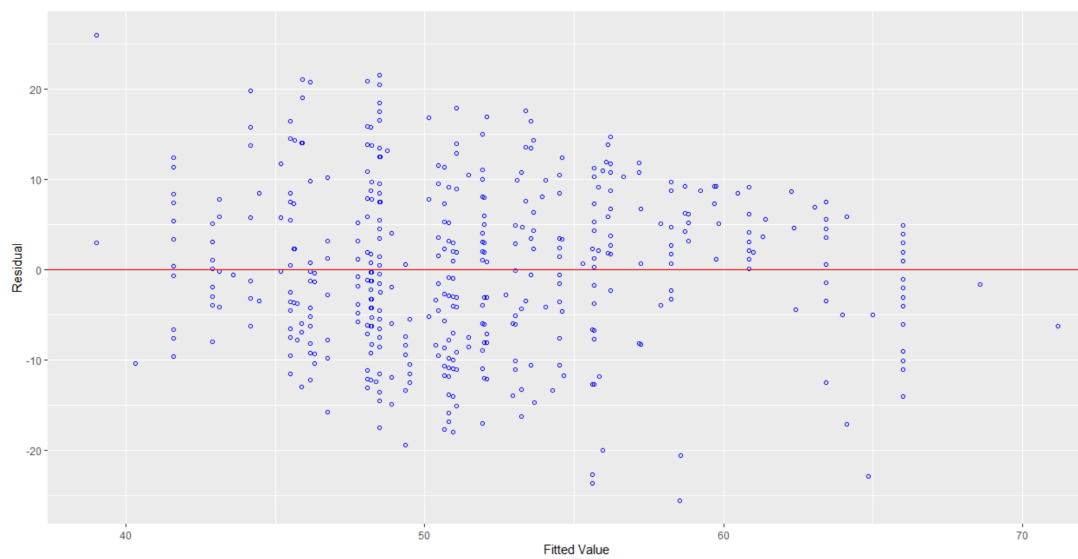


Illustration 2 – Résidus en fonction des valeurs prédites pour le modèle (3).

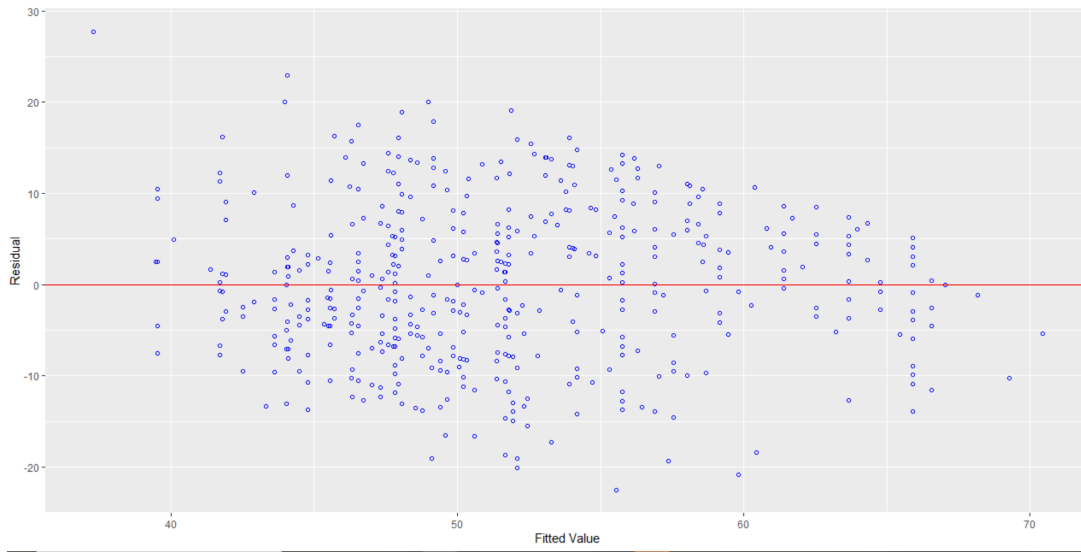


Illustration 3 – Résidus en fonction des valeurs prédites pour le modèle (10).

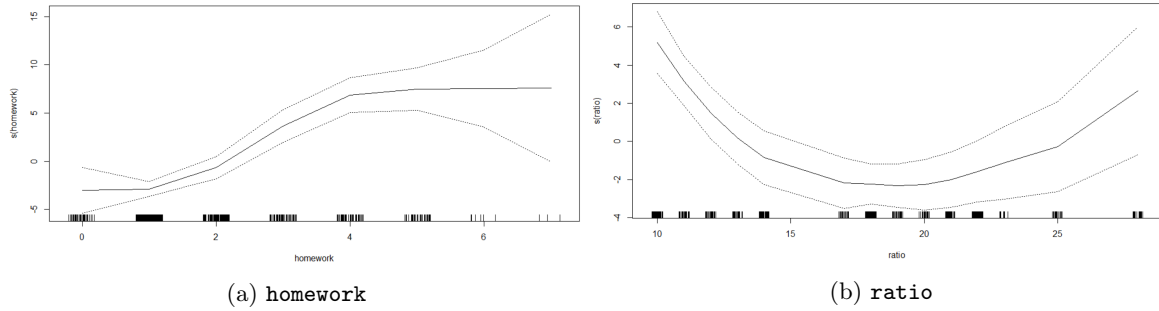


Illustration 4 – Splines réalisés sur les variables `homework` et `ratio` lors de l'entraînement d'un GAM.

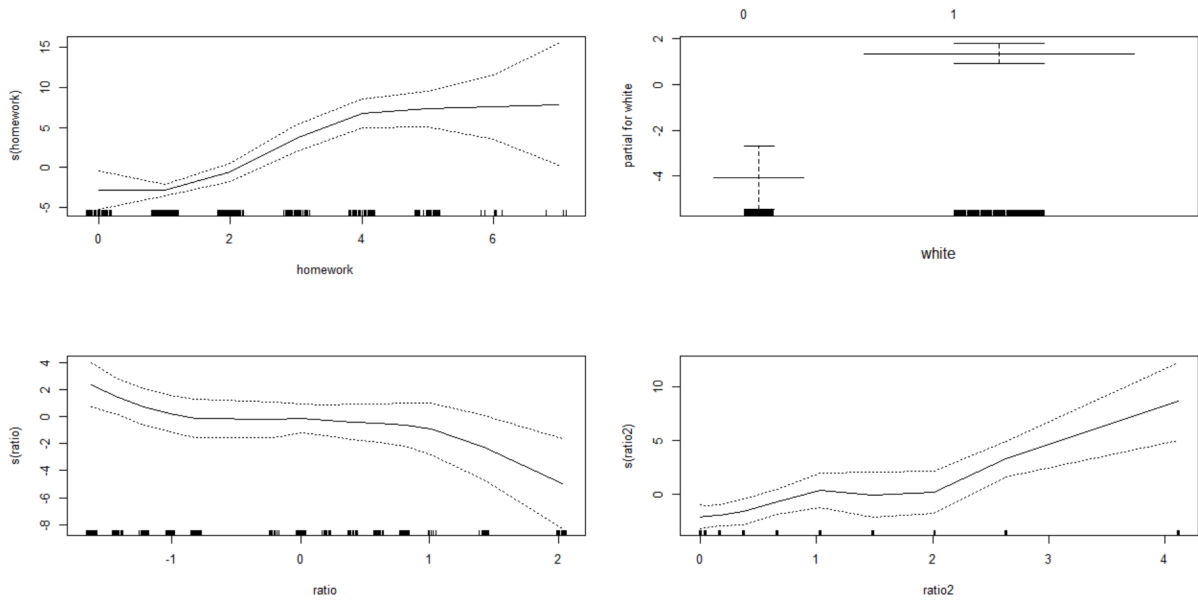


Illustration 5 – Visualisation des Splines suite à l'entraînement d'un GAM utilisant (2).

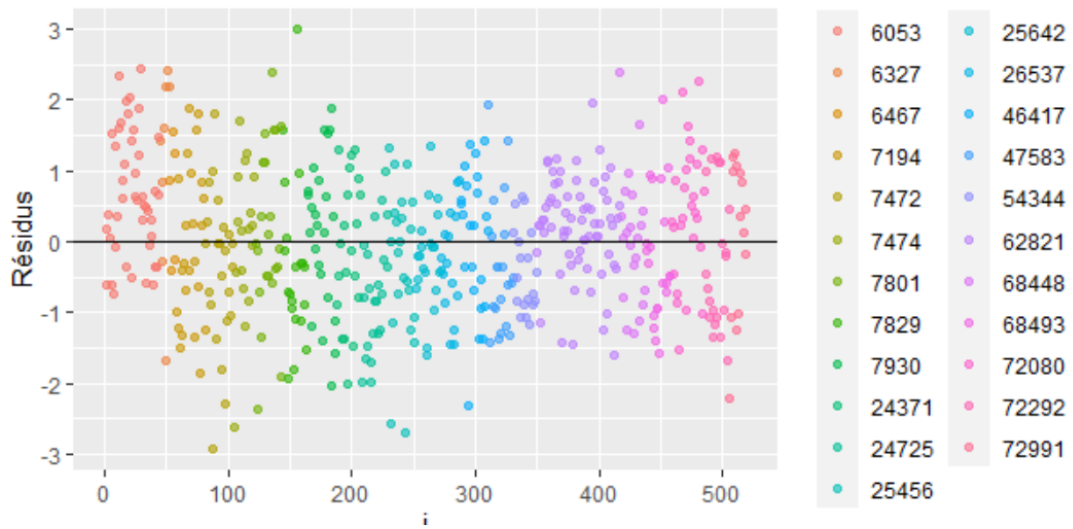


Illustration 6 – Résidus studentisés en fonction de l'index des observations pour le modèle 3.

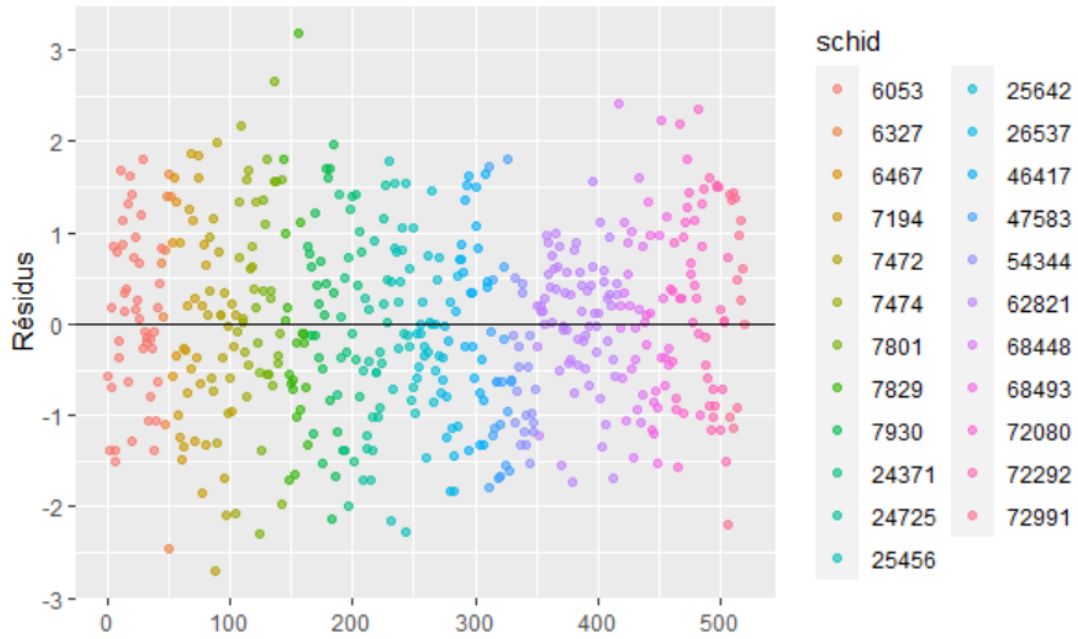


Illustration 7 – Résidus studentisés en fonction de l'index des observations pour le modèle 10.

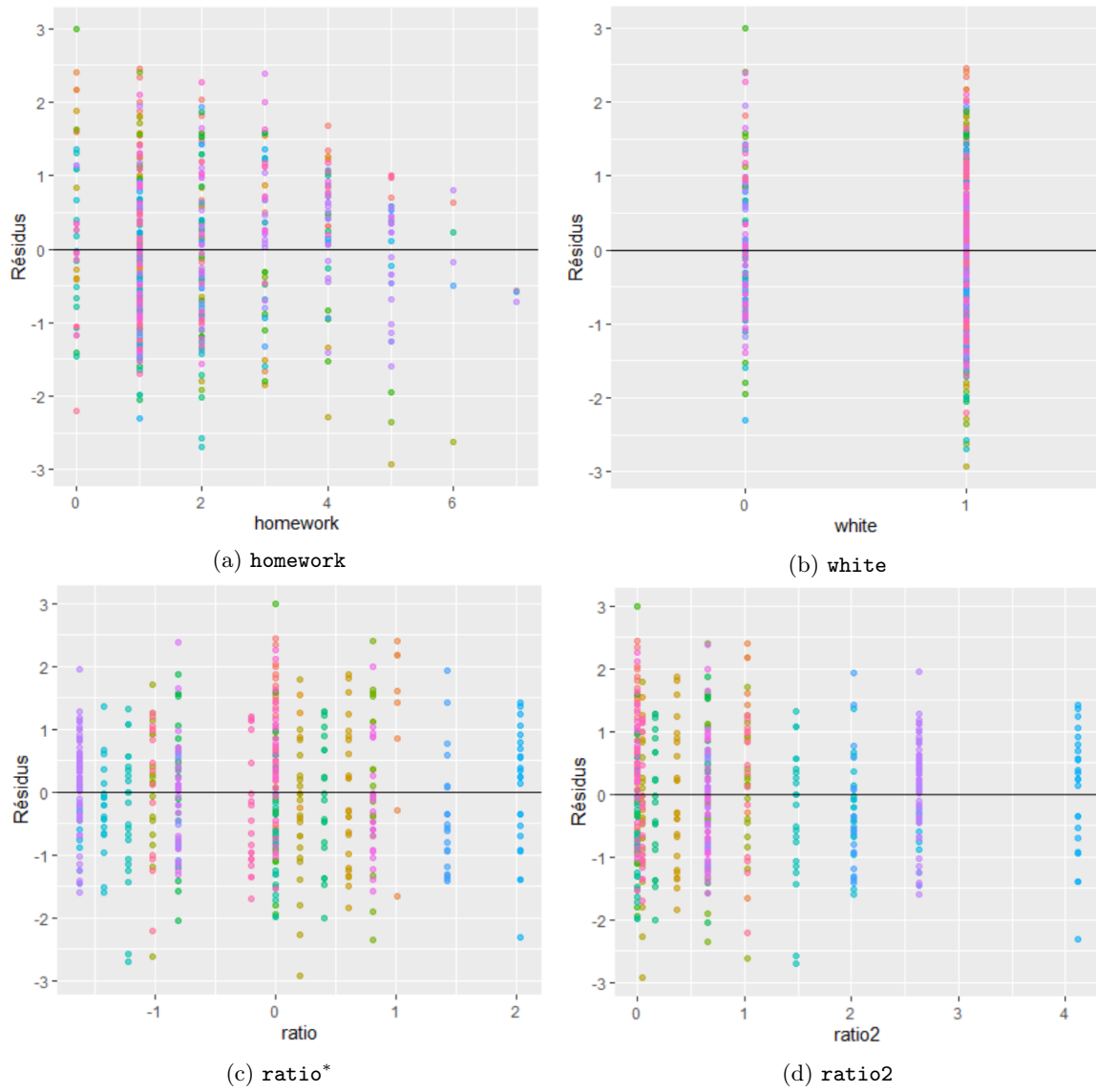
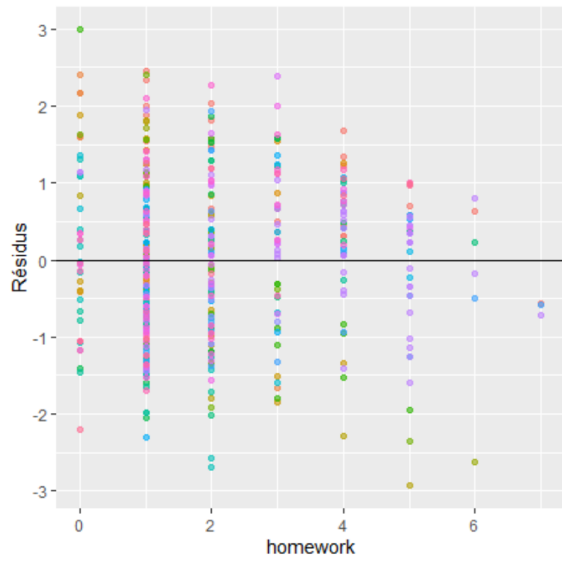
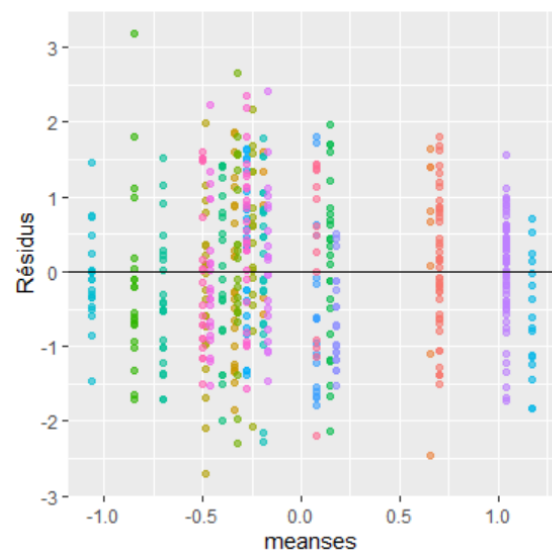


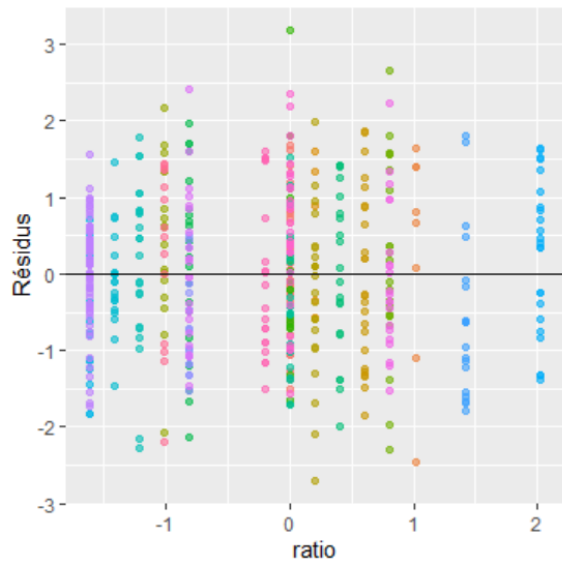
Illustration 8 – Résidus studentisés en fonction des différentes variables explicatives du modèle 3.



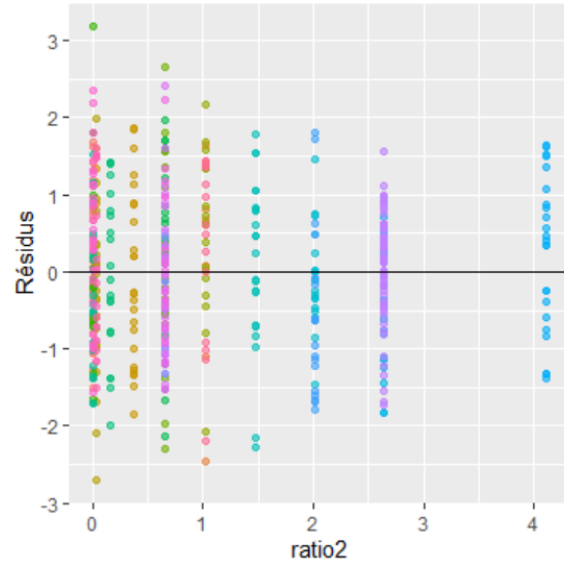
(a) homework



(b) mean ses



(c) ratio*



(d) ratio2

Illustration 9 – Résidus studentisés en fonction des différentes variables explicatives du modèle 10.


```

Linear mixed model fit by REML ['lmerMod']
Formula: math ~ homework + white + (homework | schid)
Data: data

REML criterion at convergence: 3622.8

Scaled residuals:
    Min       1Q   Median       3Q      Max
-2.30724 -0.66634 -0.03254  0.68200  3.03431

Random effects:
 Groups   Name      Variance Std.Dev. Corr
schid    (Intercept) 58.21    7.629
          homework   17.26    4.154   -0.85
Residual             52.66    7.257
Number of obs: 519, groups: schid, 23

Fixed effects:
              Estimate Std. Error t value
(Intercept)  44.0198    1.8349   23.990
homework      1.9031    0.9168    2.076
white1        3.3000    0.9781    3.374

Correlation of Fixed Effects:
          (Intr) homwrk
homework -0.773
white1   -0.371 -0.027

```

Illustration 10 – Sortie R de la fonction `summary` pour le modèle (7).

```

Linear mixed model fit by REML ['lmerMod']
Formula: math ~ meanses + homework + white + (homework | schid)

REML criterion at convergence: 3610

Scaled residuals:
    Min       1Q   Median       3Q      Max
-2.29715 -0.68843 -0.01309  0.68012  2.98973

Random effects:
 Groups   Name      Variance Std.Dev. Corr
schid    (Intercept) 53.58    7.320
          homework   16.40    4.050   -0.91
Residual             52.79    7.266
Number of obs: 519, groups: schid, 23

Fixed effects:
              Estimate Std. Error t value
(Intercept)  44.7022    1.7873   25.012
meanses       4.8925    1.3406    3.649
homework      1.9251    0.8952    2.151
white1        3.1149    0.9570    3.255

Correlation of Fixed Effects:
          (Intr) meanss homwrk
meanses   0.139
homework -0.813 -0.006
white1   -0.384 -0.126 -0.026

```

Illustration 11 – Sortie R de la fonction `summary` pour le modèle (11).

A.2 Question 2

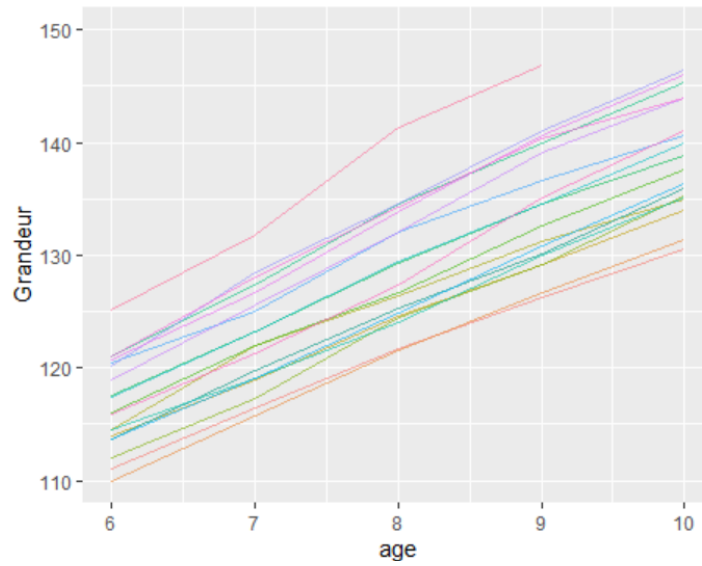
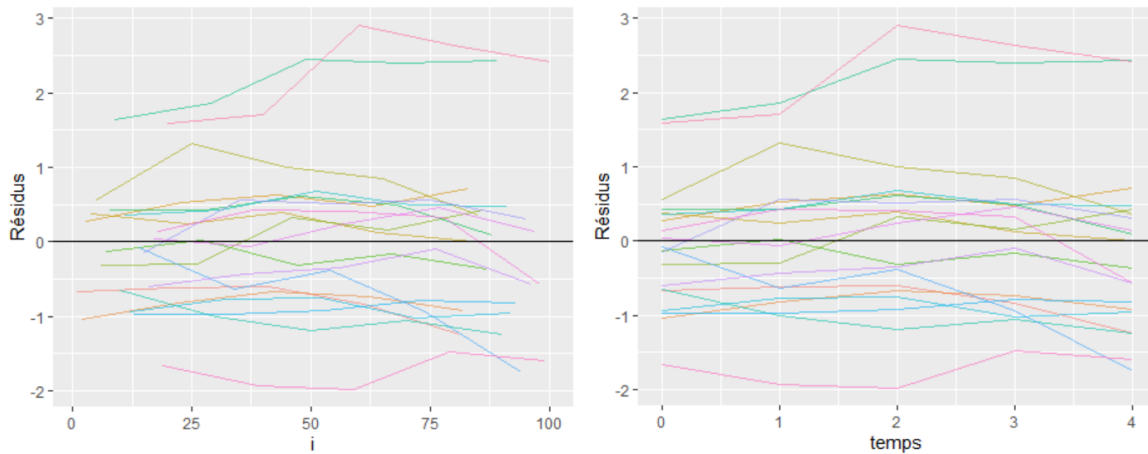
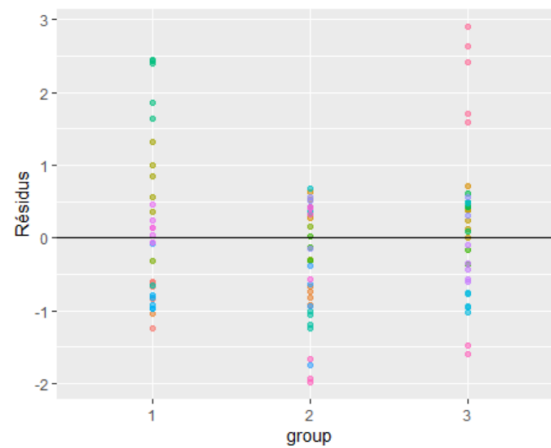


Illustration 12 – Relation de la grandeur en fonction de l'âge pour chacune des jeunes filles.



(a) i

(b) temps



(c) group

Illustration 13 – Graphiques de résidus générés à partir du modèle (15).

```

Linear mixed-effects model fit by REML
Data: data
      AIC      BIC    logLik
327.4851 350.3747 -154.7425

Random effects:
Formula: ~1 | child
      (Intercept) Residual
StdDev:  0.00225578 2.987839

Correlation Structure: AR(1)
Formula: ~1 | child
Parameter estimate(s):
      Phi
0.9498482
Fixed effects: height ~ group + temps + group:temps
              Value Std.Error DF   t-value p-value
(Intercept) 112.57455  1.2196105  77  92.30369  0.0000
group2       3.70038  1.6620490  17   2.22640  0.0398
group3       7.79495  1.6620490  17   4.68996  0.0002
temps        5.28740  0.1859984  77  28.42710  0.0000
group2:temps  0.26271  0.2534731  77   1.03643  0.3032
group3:temps  0.87029  0.2534731  77   3.43348  0.0010
Correlation:
      (Intr) group2 group3 temps  grp2:t
group2  -0.734
group3  -0.734  0.538
temps   -0.305  0.224  0.224
group2:temps  0.224 -0.305 -0.164 -0.734
group3:temps  0.224 -0.164 -0.305 -0.734  0.538

```

Illustration 14 – Sortie R de la fonction `summary` pour le modèle (17).

A.3 Question 3