

Projet Cogmont

Azat Aleksanyan, Lucas Chabeau, Matthieu François, Etienne Hamard

March 13, 2019

Contents

Remerciements	2
Introduction	3
Problématique	3
1. Contexte	4
1.1 Les besoins d'un changement éducatif	4
1.2 Les pédagogies étudiées	4
1.3 Présentation des données de départ	4
1.4 Nettoyage des données	5
1.5 Recodage des variables	5
2. Méthodologie	7
2.1 Analyse factorielle	7
2.2 Classification Ascendante Hiérarchique	7
2.3 Arbre de regression	7
2.4 Modèle linéaire mixte	7
2.5 Régression logistique avec effet mixte	8
3. Analyse exploratoire	9
3.1 Analyse Préliminaire	9
3.2 Multivariée	12
3.3 Classification Ascendante Hiérarchique des variables	17
4. Modélisation linéaire des résultats à chaque question.	20
4.1 Modélisation des questions à résultat continu par modèle linéaire mixte	20
4.2 Modélisation des questions à résultat binaire par regression logistique mixte	20
4.3 Forêt d'arbres décisionnels	21
Conclusion	22

Remerciements

Nous tenons à remercier particulièrement Madame Adeline Leclercq-Samson et Madame Marie-Caroline Croset pour leur aide, leur disponibilité et leur écoute tout au long du projet. Nous remercions aussi le CNRS pour la confiance qu'ils nous ont accordé avec ce projet.

Introduction

Le programme de notre 1ère année de Master prévoit un projet tutoré faisant appel à nos compétences acquises en statistique et en sciences des données. C'est dans ce contexte que nous avons travaillé en collaboration avec Marie-Caroline Croset et Adeline Leclercq-Samson sur un sujet mêlant sciences cognitives, statistique et fouille de données.

Depuis 4 ans, une équipe de chercheurs de l'institut des sciences cognitives (ISC) de Lyon ont testé les capacités en mathématiques des élèves d'une école maternelle située dans une zone REP+ (Réseau d'éducation prioritaire) de l'agglomération Lyonnaise. Une partie de ces élèves ont suivi une méthode d'éducation conventionnelle et les autres ont suivi des enseignements suivant la méthode Montessori. Notre objectif est de voir s'il y a ou non une différence entre le niveau en mathématiques des élèves ayant suivi les enseignements conventionnels et ceux ayant suivi les enseignements Montessori.

Nous avons reçu les résultats des tests sous forme de tableur que nous avons triés pour ne garder que les résultats des élèves de moyenne section. Nous avons ensuite nettoyé nos données et commencé l'étude sur cette population.

Problématique

Nous allons donc chercher s'il existe ou non une différence de niveau en mathématiques entre des élèves d'écoles Montessorienne et des élèves d'écoles Conventionnelle.

1. Contexte

Ce projet nous a été proposé par l’Institut des sciences cognitives - Marc Jeannerod spécialisé en neurosciences. L’UMR 5304 créée en 2007 est l’un des deux laboratoires de l’Institut des Sciences Cognitives – Marc Jeannerod. L’UMR 5304 est un laboratoire interdisciplinaire qui intègre l’expertise de chercheur des Sciences de la Vie (psychologie cognitive, neurosciences) et de médecine (pédo-psychiatrie, neuro-pédiatrie) avec celle de chercheur des Sciences Humaines et Sociales (linguistique computationnelle et théorique et philosophie) pour étudier la nature et la spécificité de l’esprit humain.

1.1 Les besoins d’un changement éducatif

Le député et mathématicien Cédric Villani a publié un rapport pour renforcer l’apprentissage des mathématiques à l’école. Les élèves français ont aujourd’hui un niveau insatisfaisant dans cette discipline. Pourtant, jusqu’en 1985, l’enseignement des mathématiques en France était reconnu comme l’un des meilleurs. L’étude internationale “Trends in International Mathematics and Science Study” (TIMSS) 2015 qui mesure les performances en mathématiques et en sciences des élèves en fin de CM1 classe la France dernière des pays de l’Union européenne. Elle obtient même un score en dessous de la moyenne internationale. Pour mettre un terme à cette tendance inquiétante de la dégradation du niveau des élèves français en mathématiques, le gouvernement est maintenant ouvert à de nouvelles pédagogie d’enseignement des mathématiques.

1.2 Les pédagogies étudiées

Pédagogie Montessori : La pédagogie Montessori est une méthode d’éducation créée en 1907 par Maria Montessori.

La pédagogie se base sur quelques principes :

-**La liberté** : les enfants sont libres de choisir l’activité qu’ils souhaitent faire parmi celles qui leur sont proposées.

-**L’autodiscipline** : les enfants sont invités à repérer eux-même leurs erreurs.

-**L’action en périphérique** : les professeurs vont préférer agir sur l’environnement de l’enfant plutôt que directement sur lui. Ils vont chercher à inciter l’enfant à faire une activité plutôt que de directement lui demander de faire cette activité.

-**Le respect du rythme de chacun** : Le rythme de l’enfant est respecté tant qu’il est concentré.

-**L’apprentissage par l’expérience** : Favoriser la pratique pour s’approprier un concept.

-**L’activité individuelle** : La plupart des activités se font en individuel.

-**L’éducation, une aide à la vie** : L’éducation est faite pour préparer l’enfant à une vie dans une société harmonieuse basée sur le respect de l’autre.

Pédagogie traditionnelle : La pédagogie traditionnelle (celle que nous connaissons aujourd’hui dans nos écoles publiques) est celle du modèle transmissif. Selon le triangle pédagogique de Jean Houssaye, cette pédagogie privilégie la relation entre l’enseignant et le savoir. Autrement dit, l’enseignant expose un savoir sous forme de cours magistral, généralement suivi d’exercices ou/et de leçons à apprendre. L’élève doit intégrer et appliquer le savoir exposé par l’enseignant.

1.3 Présentation des données de départ

Notre jeu de données est composé de trois fichiers Excel (.xlsx), avec les résultats de chaque promotion au test cognitif mis en place par l’équipe de recherche l’institut des sciences cognitives. Un quatrième fichier du même acabit pour l’année 2018/2019 nous est parvenu au milieu de l’étude.

MathsJetons__2015-2016.xlsx : Pour l’année 2015/2016.

MathsJetons_2015-2016.xlsx : Pour l'année 2016/2017.

MathsJetons_2016-2017.xlsx : Pour l'année 2017/2018.

Jetons2019.xlsx : Pour l'année 2018/2019.

Chaque jeu de données représente les résultats questions par questions (en comptant les sous-questions) des élèves ainsi que leurs catégories pédagogiques et des informations telles que l'encadrant, le niveau scolaire, la langue natale, l'âge, le type de classe (mélangé entre plusieurs section ou pas), l'année de passage du test et leur école. Il y a 10 questions divisées en sous questions, ce qui fait un total de 34 réponses. Chaque question est indépendante et pour répondre à la sous question suivante il faut une bonne réponse à la sous-question précédente, sauf pour la question 4, toutes ses sous-questions sont indépendantes. Une bonne réponse correspond à un 1 et une mauvaise réponse à un 0, sauf pour la réponse à la question 1 qui est la valeur de comptage maximale de l'enfant. Ici les élèves viennent tous d'une école située en REP+.

1.4 Nettoyage des données

Les données ayant déjà été travaillées lors d'un précédent stage, le travail de nettoyage nécessaire n'a pas été excessif. Il nous a fallu tout de même renommer certaines variables pour les rendre plus lisibles et cohérentes entre elles, supprimer certaines questions car elles n'avaient été posées qu'à certaines classes. . . Les questions étant posées de manière à ce qu'au sein d'une même tâche, il faille réussir les questions dans l'ordre pour passer à celles plus dures nous avons beaucoup de valeurs manquantes dès qu'une question était un peu difficile. Nous avons donc décidé de changer ces valeurs manquantes et les considérer comme une question que l'élève n'aurait pas réussie (donc remplacer NA par 0). Et pour ne pas passer à côté de l'information : "n'a pu aller plus loin", nous avons créé deux jeux de données : un vectoriel (chaque question devient le regroupement des résultats au sous questions. ex : $Q2.1 = 1$; $Q2.2 = 0$ devient $Q2 = 10$), un composé de scores correspondant à la somme des questions au sein d'une même tâche (pour le même exemple que le jeu de données "vectoriel" nous aurions $Q2 = 1+0 = 1$).

1.5 Recodage des variables

Afin de ne pas influencer notre jugement sur nos résultats, nous avons dans un premier temps décidé de rendre anonyme la pédagogie enseignée pour chaque classe. Chaque pédagogie fut donc renommée en "P1" et "P2". De ce fait nous n'avons pas pu privilégier une pédagogie plus qu'une autre subjectivement parlant. Aux 3 quarts du projet environ, nous avons reçu les données de nouveaux individus, une cinquantaine, et avons par la même occasion décidé d'enlever cet anonymat. Notre travail étant déjà réalisé, seul l'interprétation sur le jeu de données comportant les nouveaux individus en plus importe. Comme dit précédemment les questions sont divisées en sous questions, ces questions sont regroupables en groupe : un groupe que nous appellerons "variable au-delà", un groupe "variable outil" puis un groupe "variable objet". Chacun de ces groupes fait appelle à une tâche pédagogique en particulier.

variable Au-delà : est calculée en faisant la somme du nombre de bonnes réponses aux questions 2.3, 3.2, 4.2a, 4.2b, 4.2c, 4.2d, 5.2, 6.2, 8.6, 8.7, 8.8 et 8.9. La première question concernant la capacité à compter loin est prise en compte, les individus sachant compter au-delà de 7 ont un point en plus.

variable Outil : est calculée en faisant la somme du nombre de bonnes réponses aux questions des tâches 4, 5 et 6 soit les questions des tâches sur le surcomptage, la création d'une collection et la comparaison de deux collections.

variable Objet : est calculée en faisant la somme du nombre de bonnes réponses aux questions des tâches 1, 2, 3 et 8 soit les questions des tâches sur la capacité à savoir compter, à dénombrer une collection

et à constituer une collection d'objets. La question 1 est gérée par paliers. Nous avons séparé les individus en 5 groupes : ceux qui savent compter jusqu'à 3, puis de 4 à 7, de 8 à 10, de 11 à 16, enfin ceux qui savent compter au delà de 16. Respectivement pour chaque catégorie nous leur avons donnée un score de 0, 0.3, 0.6, 0.9 et 1.2.

2. Méthodologie

Afin de répondre au mieux à notre problématique nous avons fait le choix d'utiliser plusieurs méthodes statistiques différentes pour analyser nos données. Nous avons dans un premier temps utilisés deux méthodes qui permettent de résumer l'information globale du jeu de données : l'analyse factorielle, et la classification ascendante hiérarchique (pour faire des regroupements de variables). Puis dans un but prédictif nous avons utilisé la régression logistique mixte, la régression linéaire mixte et les forêts aléatoires. Plusieurs tests statistiques ont été réalisés en parallèle.

2.1 Analyse factorielle

Les méthodes d'analyse factorielle que nous avons utilisées ici sont l'analyse des correspondances multiples (ACM) et l'analyse en composantes principales (ACP), qui sont des méthodes de synthétisation du nombre de dimensions pour les données qualitatives et quantitatives. Cela nous permet d'appréhender plus rapidement le jeu de données, et avoir une première idée de ce qui différencie les individus entre eux (ou ce qui les rapproche). L'ACM permet dans un nuage à N dimensions, en cherchant les axes orthogonaux qui maximisent la variance entre les individus, de résumer ces N dimensions en 4 voire 5 dimensions. L'ACM a été réalisée sur les données vectorisées (regroupement des réponses aux sous-questions en un vecteur par question, question 1 exclue), celles ci ont été prises comme variables actives (celles qui définissent le placement des individus sur le graphe) et les variables portant sur la pédagogie, la question 1, et l'âge en illustratives (ajoutée après le placement des individus sur le graphe). Le principe était le même pour l'ACP qui a été faite ensuite.

2.2 Classification Ascendante Hiérarchique

N'ayant aucune information au préalable sur le thème des questions, leur regroupement... etc Mais sachant que certaines questions faisaient appelle aux mêmes compétences. Nous avons utilisé une variante de la classification ascendante hiérarchique (CAH) afin de partitionner nos variables. Nous avons aussi utilisé la CAH classique qui consiste à regrouper les individus selon leur points communs cela en partant d'une inertie interclasse maximale (inertie = somme des variances), pour arriver à une inertie interclasse de 0. Nous avons utilisé la CAH classique afin de partitionner nos individus dans un but descriptif.

2.3 Arbre de regression

L'arbre de régression est une technique d'apprentissage supervisé, qui permet en analysant un grand nombre de données, de prédire une variable à expliquer. Ils sont beaucoup utilisés dans le domaine du marketing, et plus récemment dans le domaine du machine learning (apprentissage automatique). Dans un premier temps il s'agit d'exprimer la variable à expliquer en fonction d'un maximum de variables explicatives, puis d'élaguer l'arbre afin de minimiser l'erreur, soit l'écart entre la valeur prédite et la valeur réelle. Cela revient donc à faire une régression logistique sur les données, puis d'appliquer l'algorithme de construction d'arbre à partir des résultats.

2.4 Modèle linéaire mixte

Les modèles mixtes linéaires sont une extension des modèles linéaires simples permettant des effets fixes et aléatoires. Ils sont particulièrement utilisés lorsqu'il n'y a pas d'indépendance dans les données, telle qu'elle résulte d'une structure hiérarchique.

$y_{kj} = \beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \beta_k + \varepsilon_{ij}$ sachant $k : k_{ieme}$ groupe de l'individu ; $j : j_{ieme}$ variable

2.5 Régression logistique avec effet mixte

Un GLMM (pour Generalized Linear Mixed Models) est dit “mixte”, car il comporte au moins un effet dit “fixe” (la variable dont on souhaite évaluer l’effet, ici les *Pédagogie*, *Age* et *Année Scolaire*) et au moins un effet dit “aléatoire” (la variable de regroupement, ici *newClasse* ou *Group*). Les effets aléatoires ne sont pas évalués, ils servent seulement à indiquer au modèle que les données ne sont pas indépendantes pour une boîte donnée. C’est ce qui permet à la déviance résiduelle d’être bien estimée, et ainsi à l’erreur standard des paramètres de ne pas être biaisée, et aux final d’obtenir des résultats fiables.

$$P_k(1|X) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \beta_k}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \beta_k}} \text{ sachant } k : k_{iem} \text{ groupe de l'individu}$$

3. Analyse exploratoire

3.1 Analyse Préliminaire

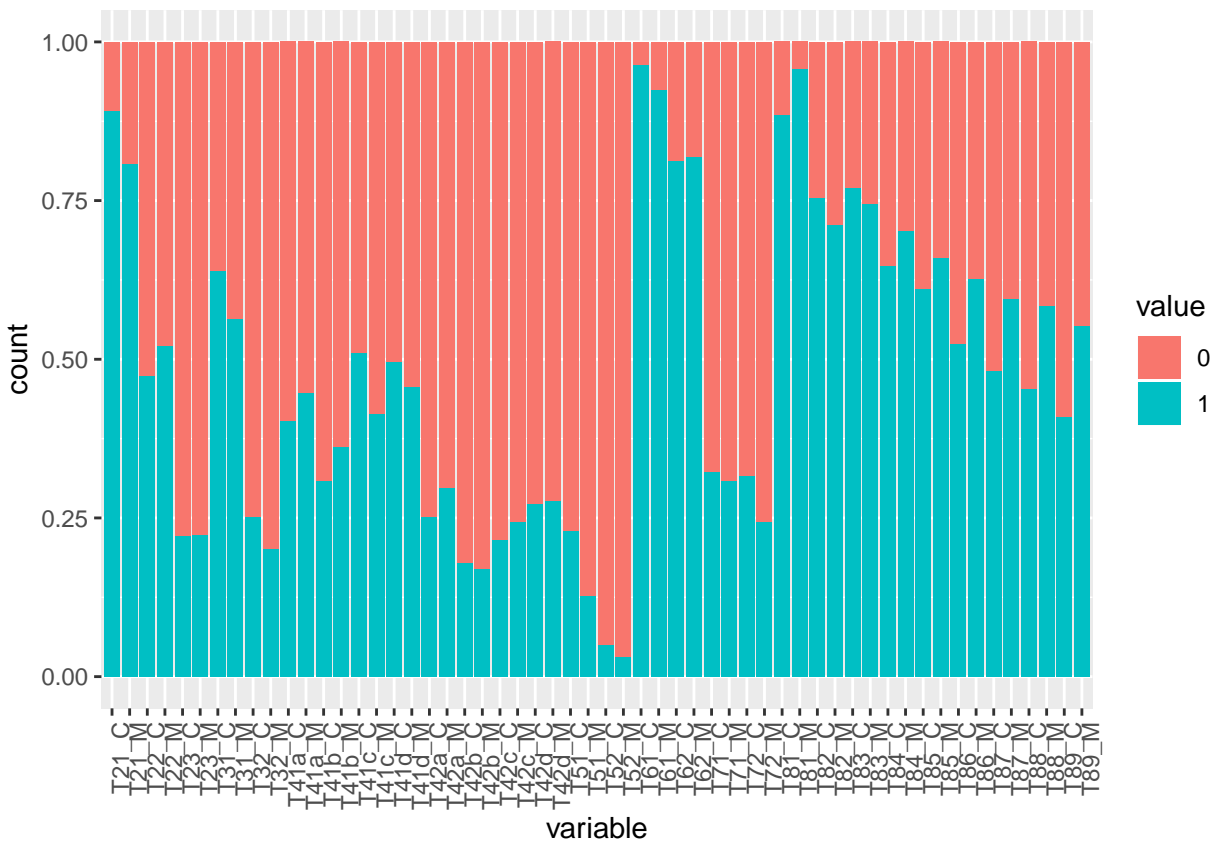


Figure 1: Réussite à chaque question par pédagogie

Le graphique précédent montre la proportion de bonnes et de mauvaises réponses par question et par pédagogie. On peut voir que l'écart entre les deux pédagogies est plus ou moins important selon les questions, mais ce n'étant pas toujours la même pédagogie qui est supérieure à l'autre, cela semble équilibré globalement. Excepté pour les dernières questions, à partir de la question 8.4, la pédagogie Montessorienne est toujours devant la Conventionnelle.

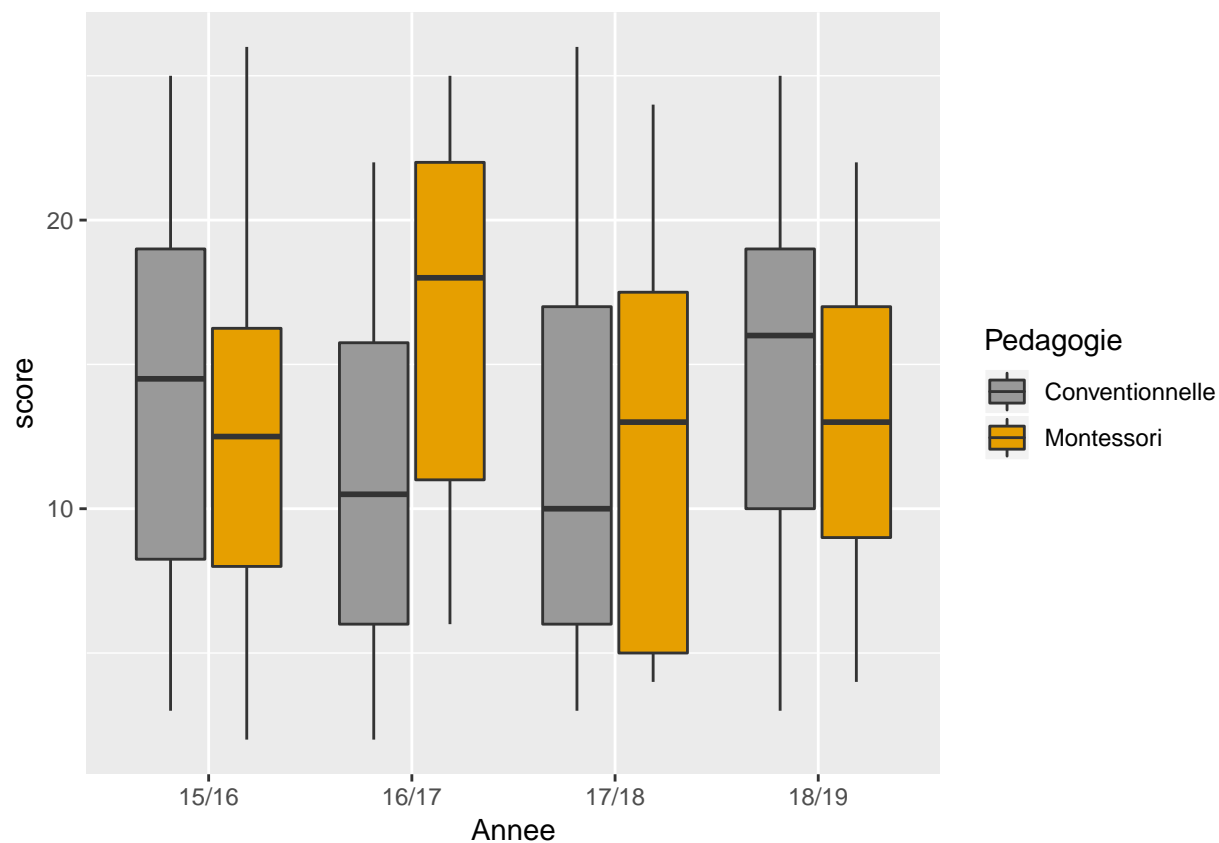
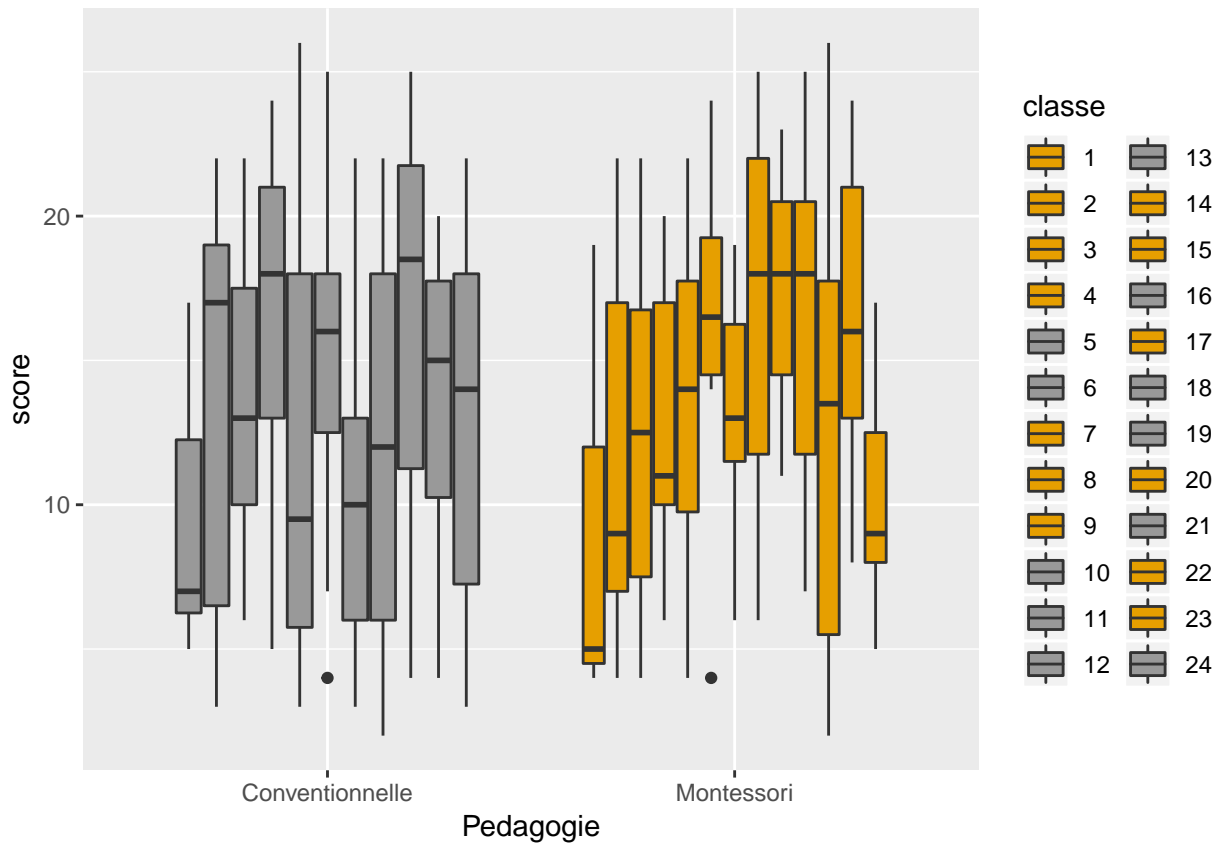
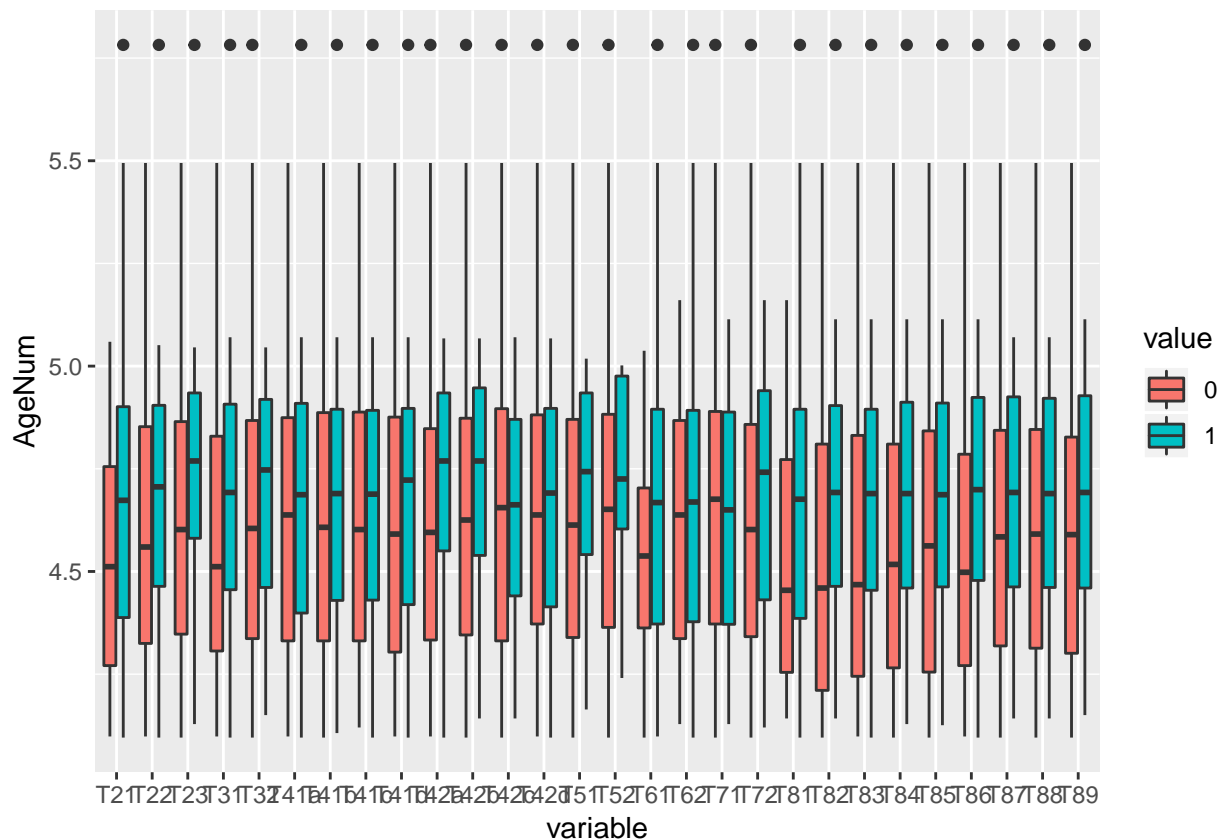


Figure 2: Score total par année par pédagogie

En ne regardant que certaines années on peut voir clairement une différence entre la pédagogie montessorienne et la pédagogie conventionnelle vis à vis du nombre de bonnes réponses. Toutefois lorsqu'on regarde la vue d'ensemble, on peut voir que selon les années, une fois la pédagogie montessorienne est supérieur à la conventionnelle, parfois c'est l'inverse. On peut donc s'attendre à ce qu'on ne puisse pas prédire quelle pédagogie permet d'obtenir un meilleur score global.



De plus il semblerait que la classe dans laquelle se trouve l'individu agisse sur le résultat de l'élève. Comme le montre le graphique précédent, la répartition des résultats totaux varie énormément d'une classe à l'autre. Enfin, après réalisation d'un test de Chi2 entre chaque question et la variable concernant le groupe des élèves, on peut observer une dépendance entre leur réussite et le groupe associé.



Enfin le graphique précédent montre que l'âge semble avoir une influence sur la réussite ou non d'une question, et ce, quelque soit la pédagogie enseignée.

Il faudra donc par la suite tenir compte de ces trois facteurs, comme des variables qui peuvent avoir un effet sur la réussite de l'élève en parallèle de la pédagogie enseignée.

3.2 Multivariée

Afin de traiter l'information présente dans le jeu de données de la meilleure façon, nous avons procédé à 2 analyses multivariées : 1 sur le jeu de données qualitatif sous forme de vecteurs, et 1 sur le jeu de données quantitatif sous forme de scores.

3.2.1 Analyse des Correspondances Multiples

La réalisation d'une ACM comme première approche sur le jeu de données a permis de mieux comprendre ce qui différencie les individus dans notre jeu de données et à la fois d'avoir un premier résultat sur la différence entre les deux pédagogies selon cette méthode.

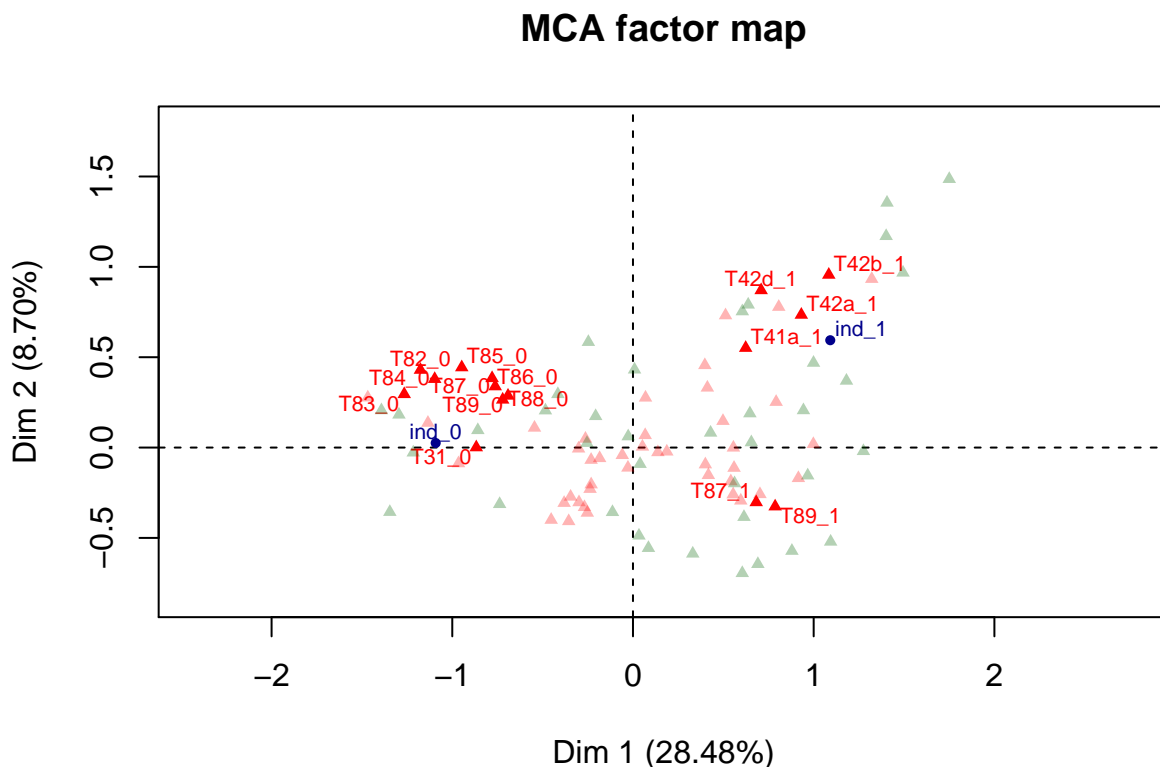


Figure 3: Graphe des modalités sur le plan principal

Le graphe précédent représente les 15 modalités qui contribuent le plus au placement des individus sur le plan principal. On peut donc voir que la dimension 1 oppose des modalités qui concernent la réussite à une question (avec un “_1” à la fin), à droite, à des modalités qui concernent l’échec d’une question (avec un “_0”), à gauche. Plus un élève réussira le questionnaire, plus il se trouvera à droite sur le graphe des individus. Cette interprétation est confirmée par l’ajout de deux individus fictifs : ind_1 et ind_0, qui comporte respectivement des succès à toutes les questions et des échecs à toutes les questions. L’individu ayant réussi en totalité le questionnaire se trouve à droite alors que l’individu ayant raté en totalité le questionnaire se trouve à gauche. De plus nous pouvons voir que les questions qui discriminent le plus la réussite ou non de l’examen sont les questions 4 et 8 (de part leur forte contribution). Toutefois cette première analyse n’aura pas permis de différencier les deux pédagogies, la variable projetée en supplémentaire sur le plan principal n’est pas significativement liée à celui ci.

Dans un second temps nous avons refait une ACM mais cette fois ci sur les données vectorielles. Cela afin de prendre en compte la succession de certaines questions qui se regroupent en “compétences”.

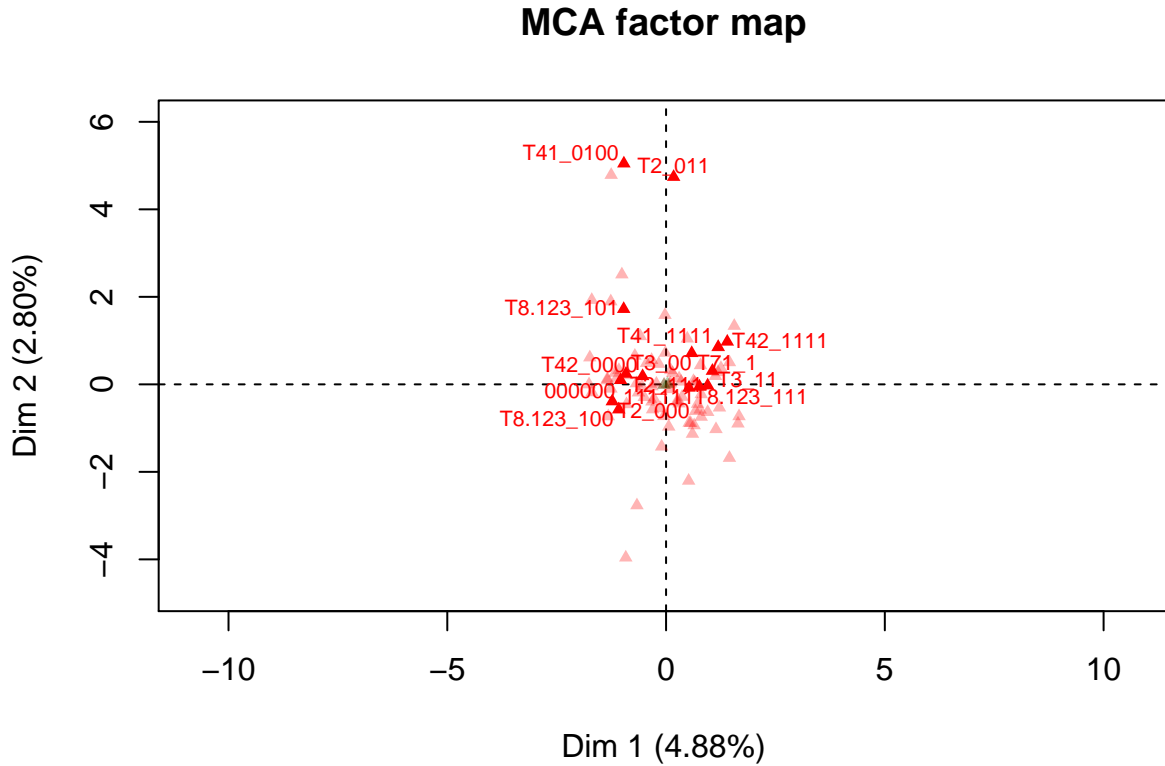


Figure 4: Graphe des modalités sur le plan principal

A nouveau la première dimension oppose les individus ayant réussi les totalités (ou la majorité pour certaines) des question à droite à ceux qui n'en ont réussi aucune à gauche. Nous ne pouvons pas non plus observer de différence significative entre les 2 pédagogies. On peut voir cette fois ci avec plus de précision les questions qui discriminent la réussite au questionnaire. Ce sont Les question 4, 3, 8 et 5.

Dans ces deux cas nous avons pu aussi observer un lien significatif entre les variables qualitatives, portant sur les réponses à la question 1 et l'âge de l'élève, et le placement des individus sur la première dimension. En conséquent on peut dire qu'il existe un lien entre la réussite à l'examen et le fait qu'un enfant sache compter "loin" et dans une moindre mesure, qu'il soit âgé.

Enfin nous avons voulu voir si en classifiant les individus suite à l'ACM nous obtenions des groupes d'individus propre à une pédagogie ou non. Pour cela nous avons utilisé la classification ascendante hiérarchique (CAH).

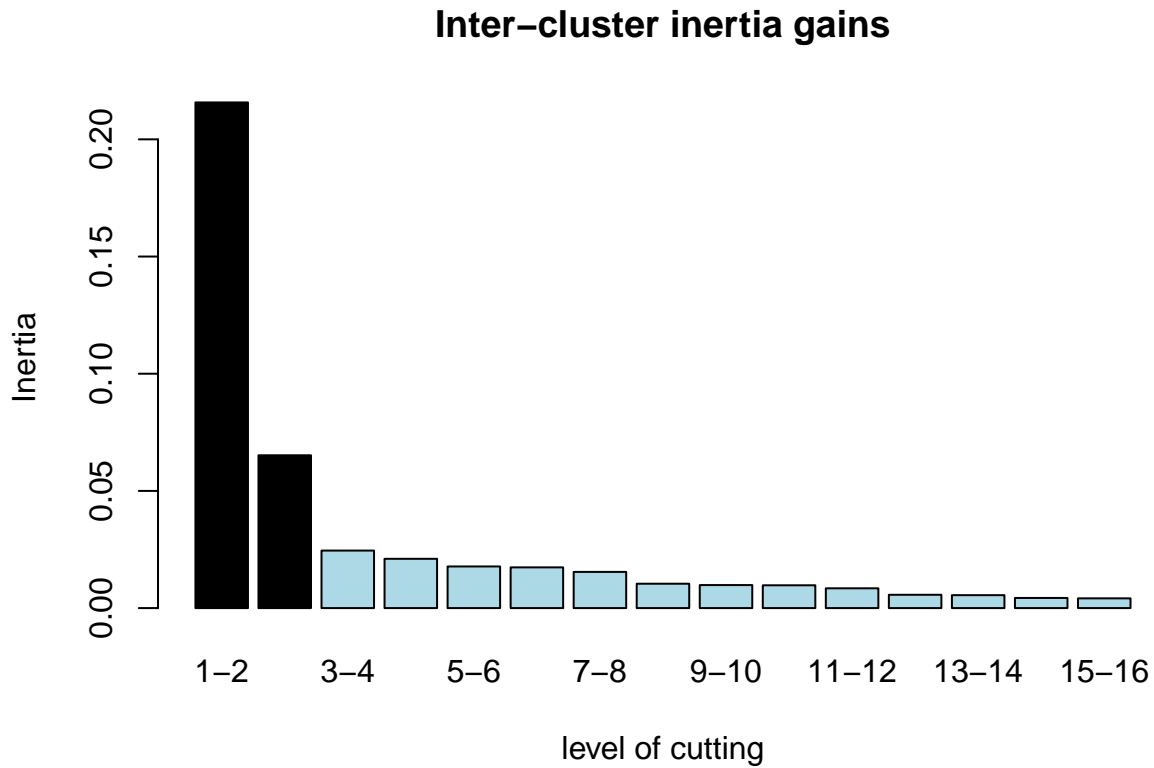


Figure 5: Diagramme des gains d’inertie

Nous pouvons observer un “saut” à la troisième classe, donc nous faisons le choix de retenir trois classes pour la CAH. Mais leur composition ne montre aucune sureprésentation d’une pédagogie plus que l’autre. Le test de χ^2 entre la variable concernant la pédagogie et celle concernant la classe n’est pas significatif. Une fois de plus cela ne permet donc pas de montrer une liaison entre la pédagogie et ce qui discrimine nos classe. Au final nous obtenions une classe d’individus qui a une majorité d’échecs, une d’individus qui échouent sur les questions 4.2, et une qui d’individus qui réussissent globalement.

3.2.2 Analyse en Composantes Principales

La réalisation d’une ACP faisant suite à l’ACM a pour but d’étudier le jeu de données différemment. En effet nous avons étudié cette fois ci le jeu de données concernant les scores. Soit un jeu de données quantitatif. Afin de ne pas perdre d’informations nous avons dans un premier temps observé la matrice des corrélations entre les variables de notre jeu de données. Car plus les variables seront corrélées entre elles, plus l’ACP ne montrera que celles ci.

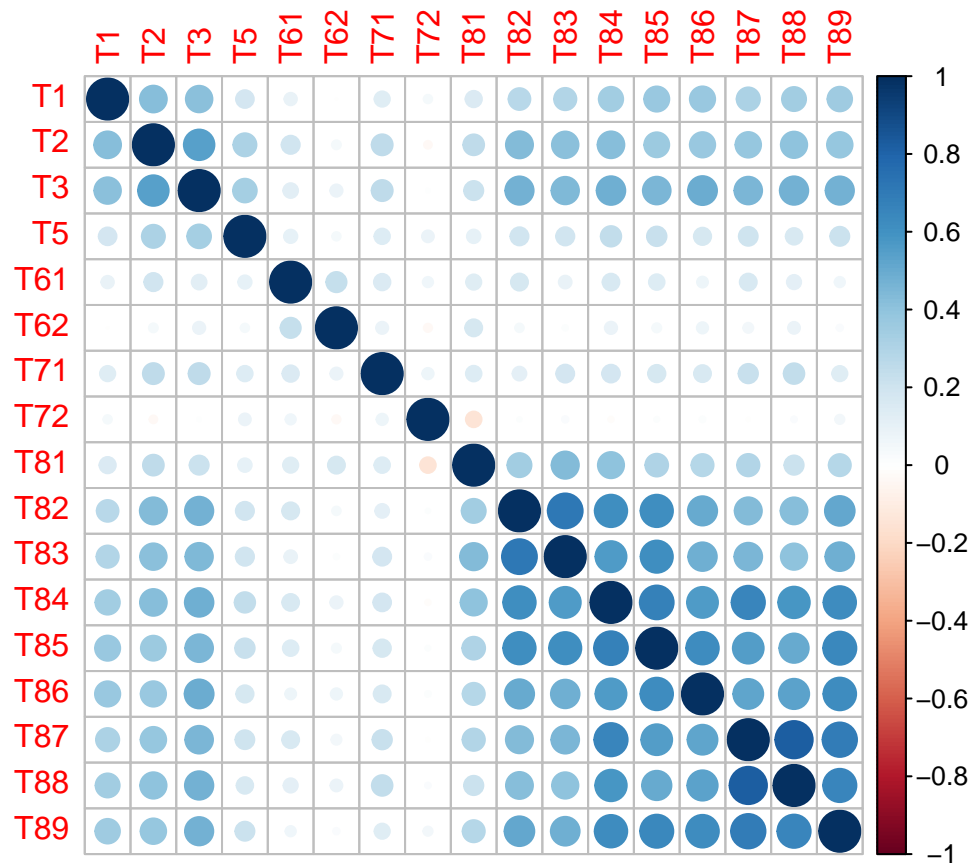


Figure 6: Matrice des corrélations

Nous avons donc fait le choix de regrouper les questions 8 en 2 groupes, aux vues des résultats. Un groupe comprenant les questions 8.1, 8.2 et 8.3, et un comportant les autres questions 8.

Variables factor map (PCA)

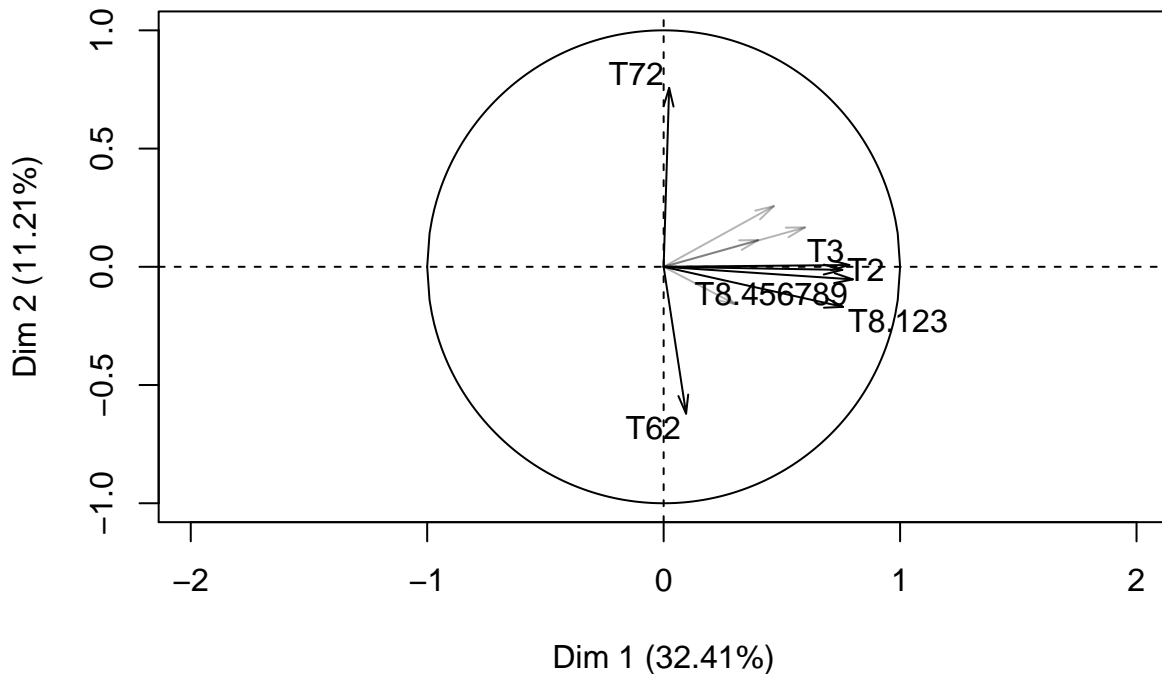


Figure 7: Cercle des corrélations du plan principal

La réalisation de l'ACP nous permet donc de voir que les individus se différencient sur la première dimension selon les questions 2, 3, et 8. Alors que la dimension 2 les différencie selon les questions 7.2 et 6.2. À nouveau, nous n'observons pas de lien significatif entre la pédagogie enseignée et le placement des individus sur les axes factoriels.

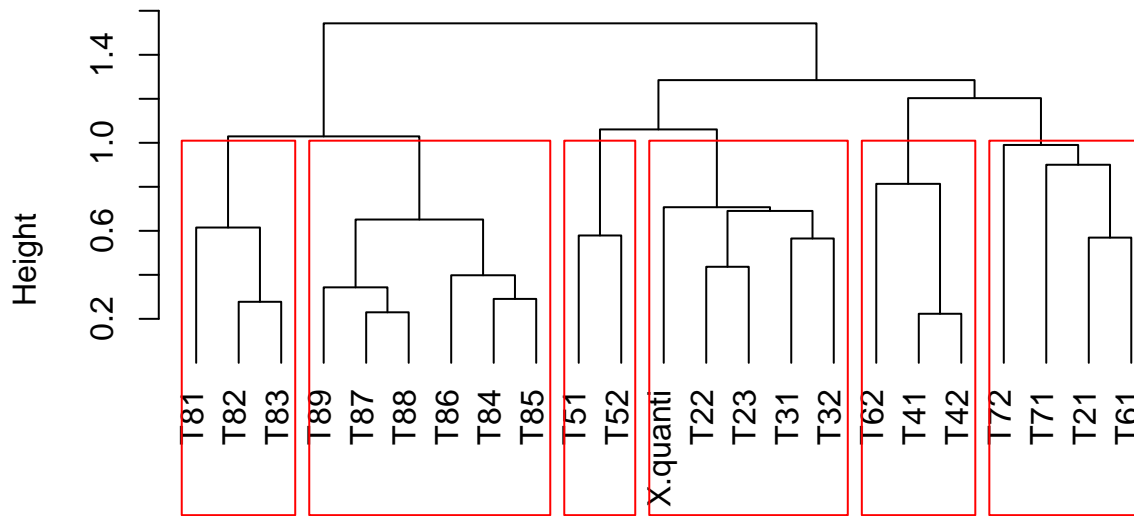
3.3 Classification Ascendante Hiérarchique des variables

N'ayant au début de notre analyse, aucune information sur le thème des questions, leur regroupement... etc. Mais sachant que certaines questions faisaient appel aux mêmes compétences. Nous avons utilisé une variante de la classification ascendante hiérarchique (CAH) afin de partitionner nos variables. La CAH est une méthode de classification qui permet de regrouper des individus dans une même classe et qu'ils soient le plus semblables possibles tandis que les classes soient elles-mêmes le plus dissemblables possibles.

Nous allons appliquer cette méthode sur les jeux de données en isolant les pédagogies pour comparer les regroupements.

Procédons d'abord à l'isolation de la Pédagogie 1.

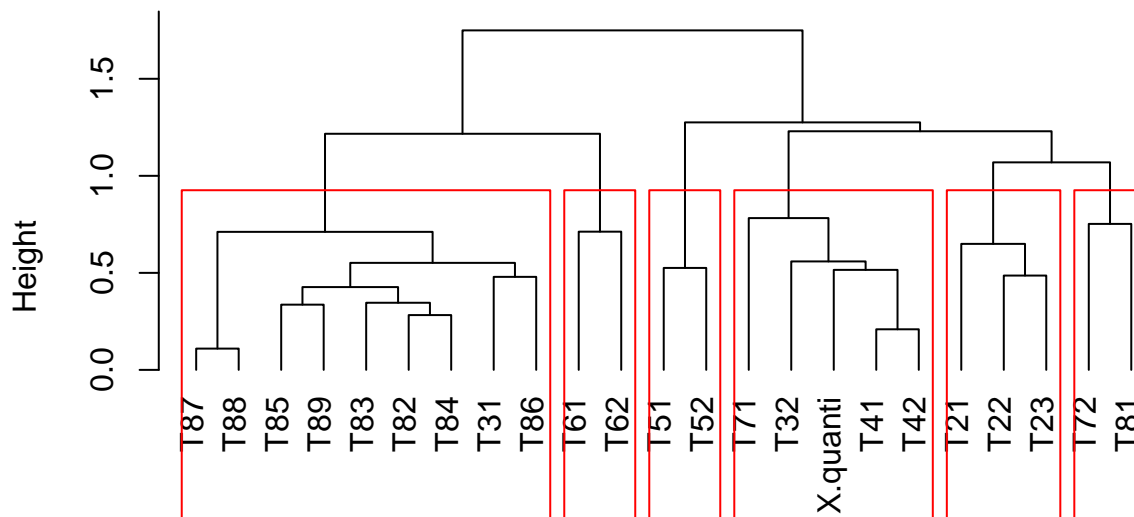
Dendrogramme des donnees generales de la pedagogie P1



Ici X.quant correspond à T1. Nous pouvons observer ici plusieurs regroupement redondant. Le regroupement des questions T81, T82 et T83 et celui des questions T84, T85, T86, T87, T88, T89. Les questions T1, T2 et T3 sont aussi fortement attirées, on retrouve en partie la variable "Objet".

Comparons maintenant avec les regroupements de la Pédagogie 2.

Dendrogramme des donnees generales de la pedagogie P2



Nous observons en regroupement des questions T82, T83, T84, T85, T86, T87, T88, T89. La T81 étant séparée du reste. On voit aussi apparaître deux couples de questions, T61 et T62 ainsi que T51 et T52. Ici nous ne voyons aucune variable prédéfinie ressortir véritablement.

On retrouve plus de similarités entre la classification de la pédagogie 2 qu'entre celles de la Pédagogie 1. Et on remarque que les regroupement qui sont stables entre les changements de jeu de données sont principalement ceux liés aux sous-questions de la T8. Pour résumer, les groupes qui ressortent sont pour la Pédagogie 1: (T81, T82, T83), (T84, T85, T86, T87, T88, T89) et (T1, T2, T3) Et pour la Pédagogie 2: (T82, T83, T84, T85, T86, T87, T88, T89), (T61, T62) et (T51, T52).

Seul la variable ‘Objet’ est en parti retrouvée et seulement dans le cas de la Pédagogie 1.

4. Modélisation linéaire des résultats à chaque question.

Nous avons voulu modéliser les résultats des élèves à chaque question en fonction de ses caractéristiques individuelles (y compris la pédagogie suivie bien-sûr) pour voir si la pédagogie a un effet sur la réussite au test. Les résultats de certaines questions sont binaires (1 pour réussite, 0 pour échec), la plupart sont concernées (toutes sauf la question 1 qui consiste à faire compter l'enfant le plus loin possible). Les résultats à l'autre question prennent des valeurs continues : la question 1 est dans ce cas mais aussi les résultats que nous avons construits en regroupant des questions (au-delà, objet et outils qui sont expliqués en 1.5).

Intuitivement, nous sommes d'abord tentés de modéliser les résultats binaires par une régression logistique et les résultats continus par une régression linéaire. Mais nous soupçonnons un effet aléatoire de la variable classe sur la réponse des élèves, c'est à dire que selon la classe (pas la pédagogie suivie) dans laquelle se trouve l'élève sa réponse sera différente. Cet effet semble être aléatoire (elle n'est pas modélisable linéairement). Il existe justement des méthodes statistiques spécialement faites pour ces cas : il s'agit des **modélisations mixtes**. Ces dernières prennent bien en compte les effets aléatoires tels la différence de résultats selon la classe de nos élèves.

Nous avons donc réalisé nos modélisations avec des **modèles linéaires mixtes** (pour les résultats continus) et des **regression logistique mixte** (pour les résultats binaires).

4.1 Modélisation des questions à résultat continu par modèle linéaire mixte

L'analyse exploratoire nous laisse en effet bien penser que les scores des élèves sont répartis différemment selon la classe dans laquelle il étudie, sa promotion et son âge. Mais seule la classe a été retenue comme effet aléatoire. L'âge a lui un effet fixe, la promotion et la pédagogie de l'élève n'ont finalement pas d'effet significatif (il s'agirait d'un effet fixe s'il était significatif).

Nous avons donc réalisé nos régressions linéaires mixtes sur les regroupements de variable en mettant l'âge comme effet de groupe (la classe) comme effet aléatoire. Et les variables sur la pédagogie, l'âge et la promotion et comme effets fixes.

Pour la modélisation de chaque regroupement de questions (au-delà, objet, outils), nous obtenons bien la présence d'un effet aléatoire, mais seule l'âge a un effet fixe sur nos variables à expliquer.

Enfin nous avons rajouté un dernier modèle sur la question concernant la capacité à compter loin. L'effet aléatoire est important, mais est justifié par l'ordre de grandeur de la variable en question.

$$\begin{aligned} Au - delà &= -6.19 + 2.493 * age + 0.79 && (effet\ groupe) \\ Objet &= -2.88 + 1.90 * age + 0.0003 && (effet\ groupe) \\ Outils &= -0.594 + 0.741 * age + 0.16 && (effet\ groupe) \\ Q1 &= -9.554 + 4.667 * age + 1.433 && (effet\ groupe) \end{aligned}$$

Pour résumer, nous n'avons pas pu modéliser la réussite à certains regroupements de tâches en fonction de la pédagogie enseignée. Toutefois nous observons qu'il est possible de les modéliser en fonction de l'âge de l'élève en prenant l'effet groupe en compte.

4.2 Modélisation des questions à résultat binaire par regression logistique mixte

N'ayant pas eu de résultats concernant la pédagogie lors de la modélisation linéaire mixte précédente, nous avons affiné notre modélisation en réalisant une regression logistique mixte sur les questions à résultat binaire.

Les principaux résultats coïncident en partie avec les résultats précédents. La pédagogie n'est significativement liée qu'à 2 questions, la question Q5.1 (création d'une collection équipotente), qui a un effet aléatoire

groupe significatif, et la Q8.9 (reconnaissance d'un chiffre particulier. **Attention !** Les questions 8.1 à 8.8 portaient sur la même thématique mais avec un autre chiffre), en complément d'un effet fixe de l'âge. De plus à nouveau plusieurs variables sont modélisable à partir de l'âge de chaque individu (voir annexe chaque modèle).

$$Q51 : P(1|Montessori) = \frac{e^{-1.1884 + -0.8875 + 0.3518}}{1 + e^{-1.1884 + -0.8875 + 0.3518}} \quad Q89 : P(1|Montessori) = \frac{e^{-6.358 + 0.6404 + 1.287*age + 0.03508}}{1 + e^{-6.358 + 0.6404 + 1.287*age + 0.03508}}$$

Il y a donc une prédisposition à réussir la question T89 pour chez les élèves ayant suivent la pédagogie Montessorienne. Inversement pour la T51 qui est en faveur de la pédagogie Conventionnelle.

Etant limité par les packages disponible de R, nous n'avons pas pu faire les courbe ROC de ces deux régressions logistiques mixtes.

4.3 Forêt d'arbres décisionnels

Les Forêt d'arbres décisionnels sont une méthode d'apprentissage automatique de régression et de classification basée sur la construction d'une multitude d'arbres de décision. Les forêts d'arbres décisionnels utilisent la méthode du bagging. Le bagging consiste à prendre un jeu de données D de taille n , et créer m nouveaux jeux de données D_i de taille n' en échantillonnant D uniformément et avec remise. Ensuite l'algorithme créer à partir de chaque échantillon D_i un arbre de décision pour ensuite agréger ces arbres et obtenir un modèle stable.

Pour prédire une variable quantitative l'agrégation ce fait par la moyenne : $G(x) = \frac{1}{m} \sum_{i=1}^m G_i(x)$.

Et pour prédire une variable qualitative on procède à un agrégation par le vote : $G(x) = \text{Vote majoritaire } (G_1(x), \dots, G_m(x))$ où $G_i(x)$ représente un modèle entraîné sur un ensemble D_i .

Nous avons donc appliqué cet algorithme sur deux jeux de données, le jeu de données contenant les variables regroupées et le jeu de données de base contenant toutes les questions une à une.

Nous avons ici la matrice de confusion. Cette matrice nous donne en ligne les données observées et en colonnes les données prédites. Il y a 64 individus issus de la pédagogie Conventionnelle ayant été bien classés et 32 individus ayant été mal classés et considérés comme des Montessori. Ce qui nous donne un taux de bonne prédiction pour les individus de la pédagogie Conventionnelle de 33%.

Pour la pédagogie Montessori, 26 ont été bien classés et 41 ont été classé en Conventionnelle à tort. On a donc une minorité de bonne prédiction. Nous ne pouvons pas obtenir de prédiction efficace grâce à ce modèle.

##	Conventionnelle	Montessori	class.error
## Conventionnelle	64	32	0.3333333
## Montessori	42	25	0.6268657

Maintenant voyons pour le jeu de données avec les questions générales.

La matrice de confusion nous indique qu'il y a 24% de mauvaise prédiction pour la catégorie Conventionnelle et 76% de mauvaise prédiction pour la catégorie Montessori

##	Conventionnelle	Montessori	class.error
## Conventionnelle	80	20	0.2000000
## Montessori	50	13	0.7936508

Conclusion

Afin de d'étudier les éventuelles différences ou non en mathématiques entre les élèves suivants une pédagogie Montessorienne et et ceux suivants Conventionnelle, nous avons commencé par une brève étude exploratoire permettant de comprendre nos données et de repérer les possibles liens entre les variables. Par la suite une analyse prédictive sur les réponses aux questions en fonction de la pédagogie aurait pu permettre de différencier ces deux pédagogies. Toutefois, les résultats de cette analyse sont peu concluants. Seuls 2 modèles différencient les 2 pédagogies : la modélisation des réponses aux questions concernant la création d'une collection équipotente et la reconnaissance d'un chiffre particulier. Il serait intéressant de continuer cette étude avec un nombre d'individus plus important voir si de nouveaux liens apparaissent.