

# Report V1

Congo Job

May 27, 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Sport et sciences sociales . . . . .	1
1.2	Main Objectives . . . . .	2
1.3	Specific Objectives . . . . .	2
<b>2</b>	<b>Description of raw data</b>	<b>2</b>
<b>3</b>	<b>Preprocessing</b>	<b>3</b>
3.1	Preprocessing Method . . . . .	3
3.2	Results of processed data . . . . .	4
<b>4</b>	<b>Clustering</b>	<b>4</b>
4.1	Principal component analysis . . . . .	4
4.2	Hierarchical Clustering on Principal Components . . . . .	11
4.3	Results of HCHC . . . . .	13
<b>5</b>	<b>Classification</b>	<b>14</b>
5.1	Choice of Classification algorithm . . . . .	14
5.2	Feature selection . . . . .	16
5.3	Results final model . . . . .	17
<b>6</b>	<b>Conclusion</b>	<b>18</b>

## 1 Introduction

La pratique rgulire dactivits physiques comporte de nombreux bnfices tel qu' une amlioration de la sant mentale ,la prvention des maladies cardiovasculaire , limiter la prise de poids et bien d'autres .

Cependant, on observe un dclin de la pratique dactivits physiques. il apparat primordial dencourager les jeunes maintenir une activit physique ou devenir plus actif , c'est dans cette optique que s'inscrit ce projet .

### 1.1 Sport et sciences sociales

Cre en 1979 par Bernard Michon Strasbourg ,lunit de recherche Sport et sciences sociales demeure la seule unit de recherche STAPS du Grand Est et est reconnue comme une structure de recherche incontournable en sciences sociales du sport dans le paysage franais et europen.

Regroupant plus de 20 chercheurs (titulaires et associés) et 17 doctorants, elle réalise des ouvrages et des articles de référence (plus de 174 publications).

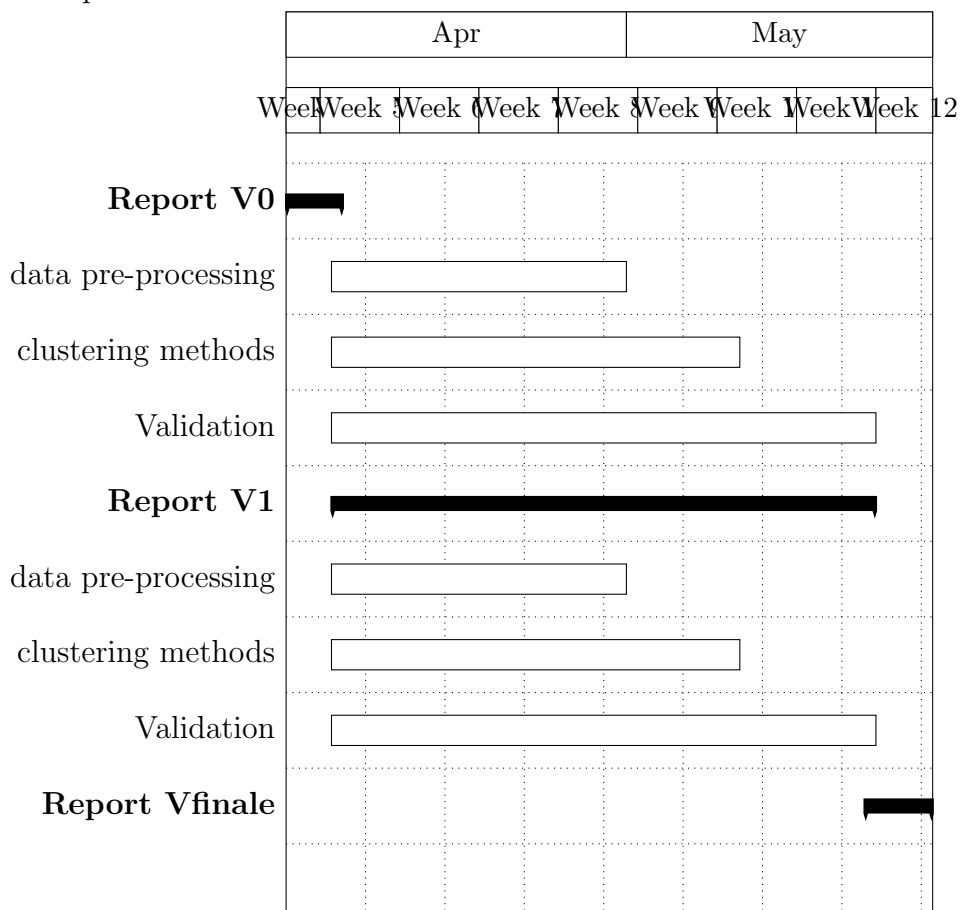
## 1.2 Main Objectives

L'objectif de ce travail est d'identifier des profils de praticiens basés sur des qualificatifs positifs ou négatifs, c'est-à-dire d'attribuer un profil à chaque cluster des données et d'estimer la force de ces profils, c'est-à-dire le nombre de clusters ou de profils qui sont les plus représentatifs des données.

## 1.3 Specific Objectives

Nous allons d'abord effectuer un prétraitement des données par la renormalisation des données, la suppression des valeurs aberrantes, compléter ou supprimer les valeurs manquantes, puis pour analyser les données nous utiliserons différents algorithmes tels que : K-means, analyse en composantes principales, arbres de décision. Enfin, nous testerons la solidité de notre cluster en utilisant des algorithmes de classification tels que : la régression logistique, les k- plus proches voisins,...

Le diagramme de Gantt ci-dessous nous donne un aperçu rapide de l'organisation du travail dans le temps.



## 2 Description of raw data

Le dataset contient des informations personnelles sur les lycéens (1070 participants) tel que leur initial, le lycée, le sexe, le choix d'étude, le travail des parents et le support des parents ainsi que la date de naissance, la morphologie de la personne (la taille et le poids). Une vingtaine

des variables mesurent la nature de la motivation par exemple la jouissance , l'affiliation , la condition physique et le dgr de motivation tel que SIMS intrinsic et SIMS external regulation .

Enfin le reste des variables (71) recolt de la manire suivante : On pose une question : "En EPS, quel est le sport que vous avez le plus apprci ?".

Puis on lui indique : "Nous allons maintenant te prsenter des mots qui vont te permettre de dcire ton ressenti par rapport ce sport. Votre travail consiste indiquer, le plus rapidement possible, si vous tes d'accord ou non avec ces propositions en cliquant sur oui ou non. Le temps de rponse a t pris en compte dans chaque rponse. Si ce temps est court, cela signifie que le terme semble vident. Par exemple, si le sport est "le football", l'lve pourra rpondre "oui" rapidement au qualificatif "plaisir", "non" rapidement au qualificatif "beaut".

Les rponses possibles chaque question sont "oui", "non", "je ne sais pas". Lorsque la rponse une question est "oui", la valeur du temps est positive, ngative dans le cas de "non" et zro dans le cas de "je ne sais pas".

	AP	AQ	AR	AS	AT	AU	AV	AW
1	Qualite	Force	Maintien	Puissance	Compétition	Muscle	Beaute	Galbant
2	2 269	2 637	1 271	1 330	1 297	1 087	1 376	6 982
3	1 621	1 135	1 426	1 444	1 134	1 329	1 394	1 329
4	2 156	2 627	2 674	3 858	2 886	2 676	2 869	9 192
5	1 083	1 316	1 134	1 199	1 640	1 531	2 084	1 916
6	1 232	2 660	2 130	1 517	1 297	1 577	1 633	3 339
7	1 176	1 337	1 329	1 073	970	2 011	1 003	990
8	1 548	1 492	1 540	1 377	1 062	817	1 007	21 098
9	1 784	954	2 818	1 200	2 575	2 405	1 476	2 951
10	846	1 910	1 897	1 459	917	2 065	924	1 394
11	114	97	48	153	26	178	35	127
12	815	681	1 018	648	606	433	210	195
13	1 995	1 523	1 442	811	1 183	1 070	4 232	8 713
14	4 630	1 188	1 298	973	939	805	1 054	1 362
15	1 005	876	2 027	3 372	3 534	1 464	968	6 096
16	13 230	-24 103	28 141	0	0	0	0	53 420
17	7 633	-9 876	12 933	14 221	15 321	-19 314	22 099	-26 416
18	6 241	7 827	9 013	10 382	0	0	16 164	17 232
19	7 378	9 059	10 652	0	14 980	16 124	0	20 313
20	0	0	0	0	0	0	0	0
21	-8 010	9 132	0	0	0	0	17 893	0
22	0	0	0	0	0	0	0	0
23	7 700	8 453	9 783	10 686	11 565	12 252	14 208	-20 362
24	1 236	0	-2 151	0	1 281	-2 718	0	0
25	1 951	1 350	1 651	1 068	1 351	1 800	0	0

Figure 1: Raw data

## 3 Preprocessing

### 3.1 Preprocessing Method

Les valeurs manquantes sont mis zros en supposant que le qualificatif n'intresse pas les tudians concerns et qu'ils auraient pu rpondre : "je ne sais pas". Pour la gestion des valeurs aberrantes , celles qui sont suprieur 5\*cart-type ont t mise zros en fin de ne pas impacter le poids donne chaque mot. L'cart-type est calcul en utilisant les donnees non signes dans le but de rduire les valeurs extrme et viter la possible compensation des valeurs. Les tudians ayant rpondu "je ne sais pas" toutes ces questions n'ont pas t considr dans la suite du projet . La normalisation est faite par ligne dans le but de conserver ce qui est "important" pour chaque personne.

## 3.2 Results of processed data

Qualite	Force	Maintien	Puissance	Compétitivité	Muscle	Beauté	Galbant
0,07968	0,09833	0,02909	0,03208	0,03041	0,01977	0,03441	0,31855
0,37034	0,25278	0,32317	0,32753	0,25254	0,29971	0,31543	0,29971
0,06264	0,08757	0,09006	0,15275	0,10129	0,09017	0,10039	0,43517
0,15351	0,21693	0,16739	0,18508	0,30512	0,27545	0,42597	0,38024
0,04507	0,15859	0,11645	0,06773	0,05024	0,0725	0,07695	0,21256
0,04796	0,06543	0,06456	0,03678	0,02561	0,13856	0,02919	0,02778
0,04451	0,04185	0,04413	0,0364	0,02147	0,00986	0,01886	0,97113
0,20936	0,04816	0,41018	0,09594	0,36298	0,32997	0,14954	0,43601
0,171	0,44751	0,44413	0,3303	0,18945	0,48779	0,19127	0,31341
0,01237	0,01039	0,00467	0,01692	0,0021	0,01984	0,00315	0,01389
0,07193	0,05717	0,09429	0,05354	0,04891	0,02985	0,00529	0,00364
0,1346	0,0999	0,09395	0,04756	0,07491	0,0666	0,29905	0,62846
0,83888	0,09546	0,11922	0,04903	0,04168	0,01274	0,06652	0,13305
0,01576	0,01031	0,05894	0,11577	0,12261	0,03515	0,0142	0,23086
0,05954	-0,10847	0,12664	0	0	0	0	0
0,03856	-0,0499	0,06534	0,07185	0,07741	-0,09758	0,11165	-0,13346
0,08063	0,10111	0,11644	0,13412	0	0	0,20882	0,22262
0,06209	0,07624	0,08965	0	0,12607	0,1357	0	0,17095
0	0	0	0	0	0	0	0
-0,03752	0,04278	0	0	0	0	0,08382	0
0	0	0	0	0	0	0	0
0,07073	0,07764	0,08986	0,09816	0,10623	0,11254	0,13051	-0,18703
0,29527	0	-0,51386	0	0,30602	-0,64931	0	0
0,70484	0,48772	0,59646	0,38584	0,48808	0,65029	0	0

Figure 2: Data processed

Les mots les plus "importants" et les moins "importants" pour chaque personne sont respectivement proche de soit 1 ou de -1 . Ceux qui sont moins "importants" sont proches de zero.

Le jeu de données nettoiyé contient 1050 lycens et 71 caractéristiques.

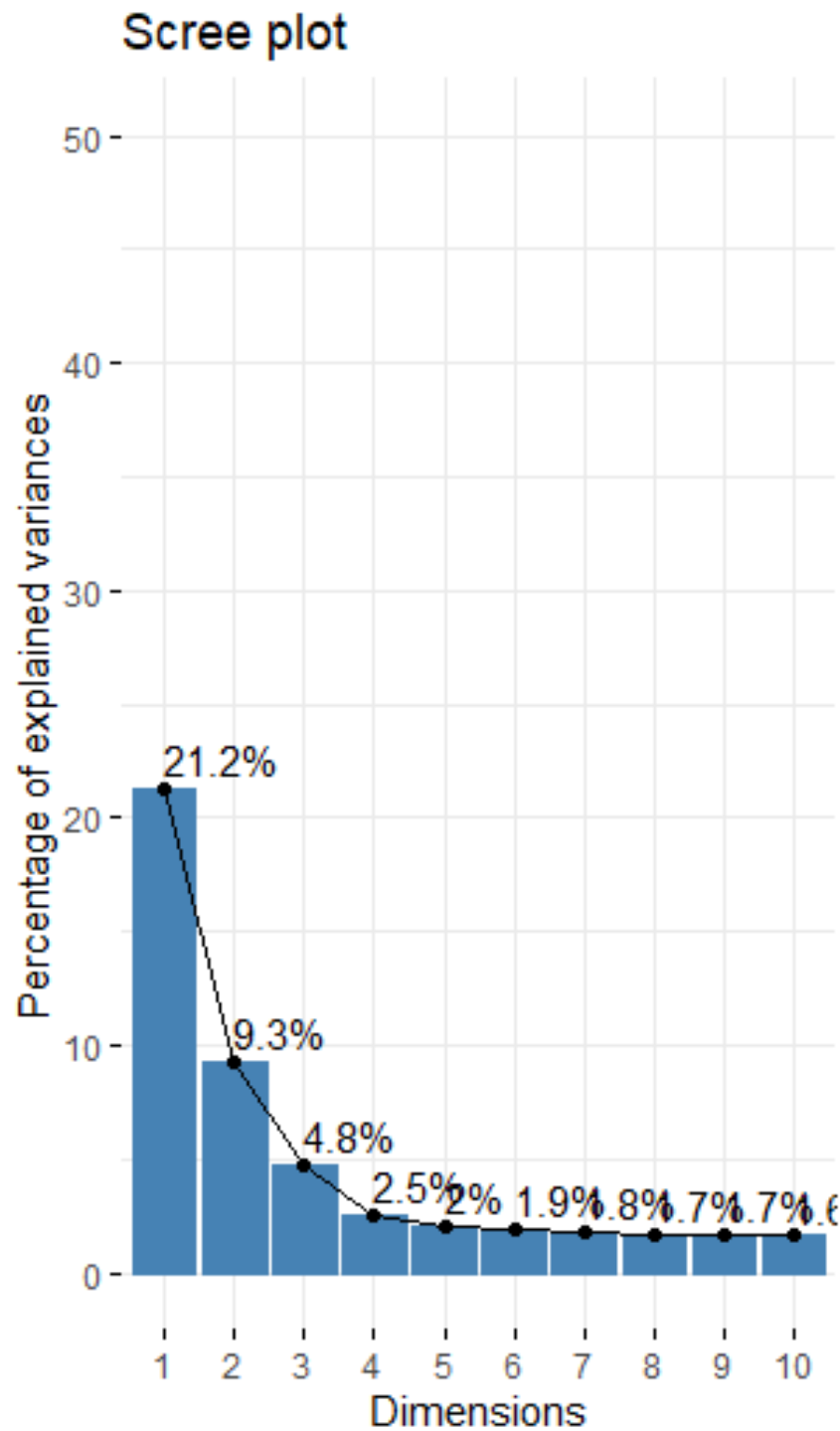
## 4 Clustering

### 4.1 Principal component analysis

L'analyse en composantes principales (ACP) est utilis pour rduire la dimension des données en quelques variables et garder les données les plus importants. C'est une mthode de la statistique multivarie, qui consiste transformer des variables lies entre elles (dites corrles en statistique) en nouvelles variables decorrles les unes des autres.

Pour dterminer le nombre de composante optimale , la comande fviz\_eig de Rstudio, a t solicit. Cette fonction permet d'avoir le graphique des valeurs propres.

Les valeurs propres mesurent la quantité de variance expliquée par chaque axe principal. Elles sont grandes pour les premiers axes et petites pour les axes suivants.



**Figure 3:**

D'après le graphique ci-dessus, nous pourrions vouloir nous arrêter à la cinquième composante principale car la variation est moindre après la cinquième. Cependant 39.79760 % des informations (variances) contenues dans les données sont retenues par les 5 premières composantes principales.

Le graphique ci-dessous montre le top 30 des variables contribuant le plus aux 5 composantes principales. Les lignes en pointillés rouges, sur les graphiques, indiquent la valeur contribution moyenne.

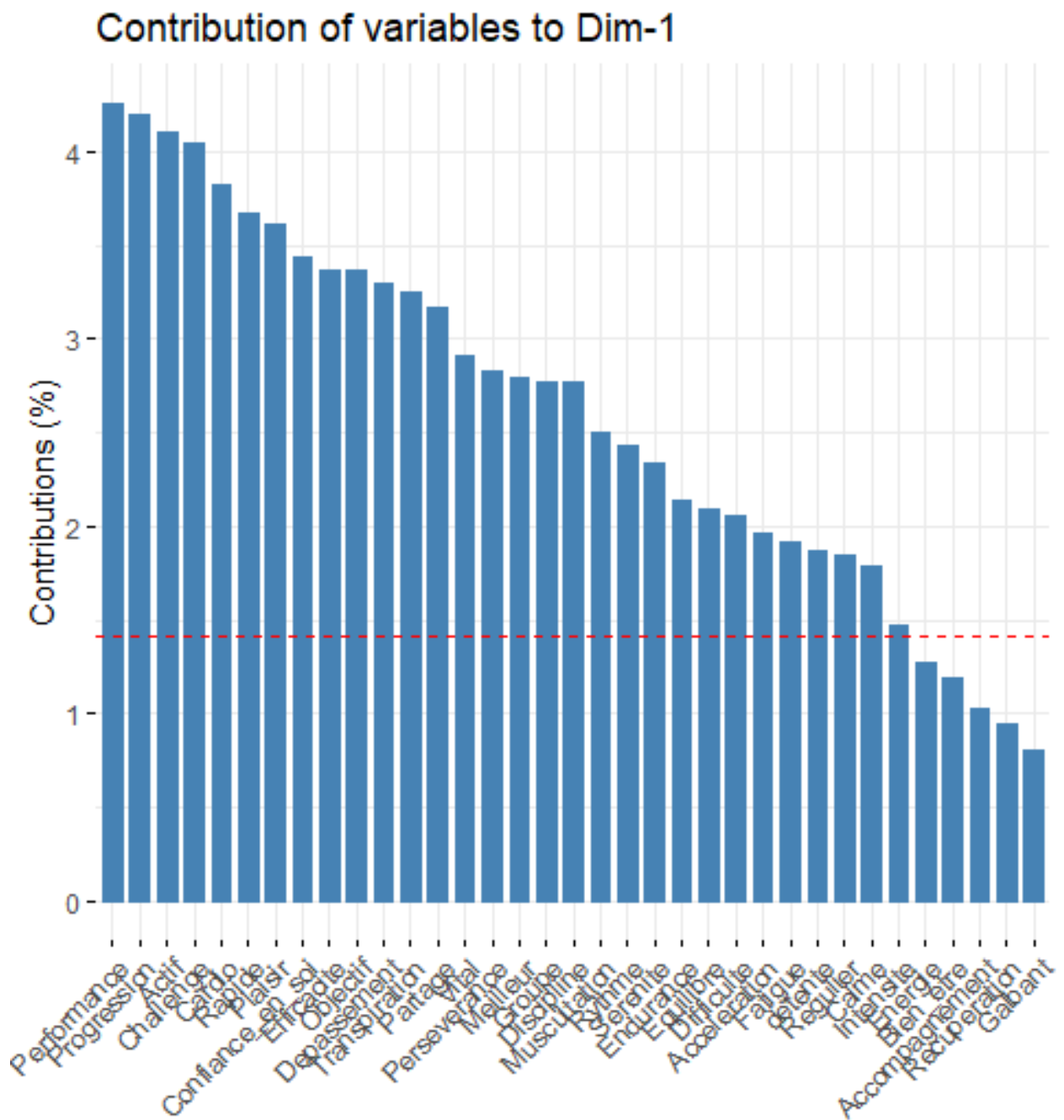


Figure 4:

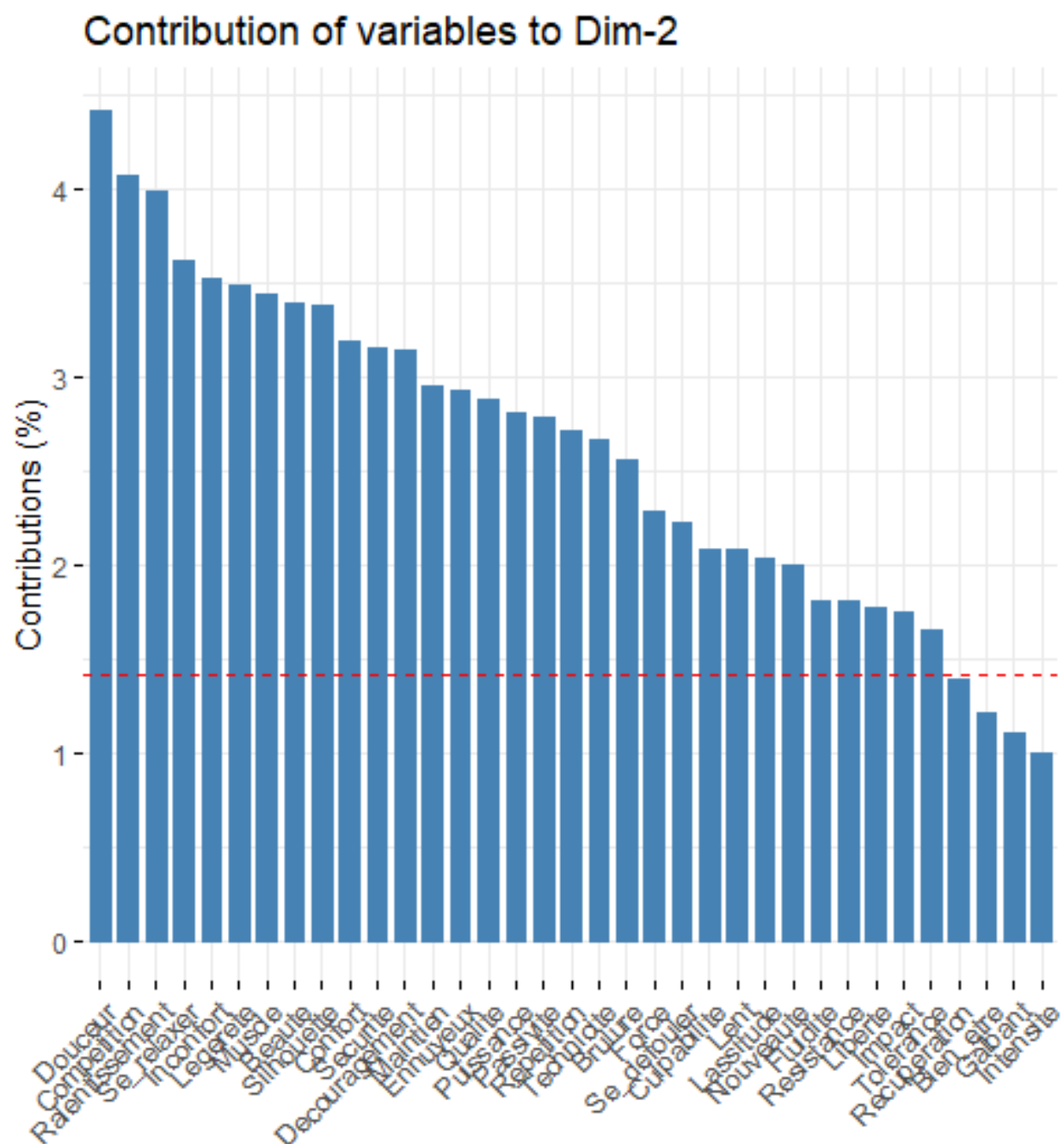


Figure 5:

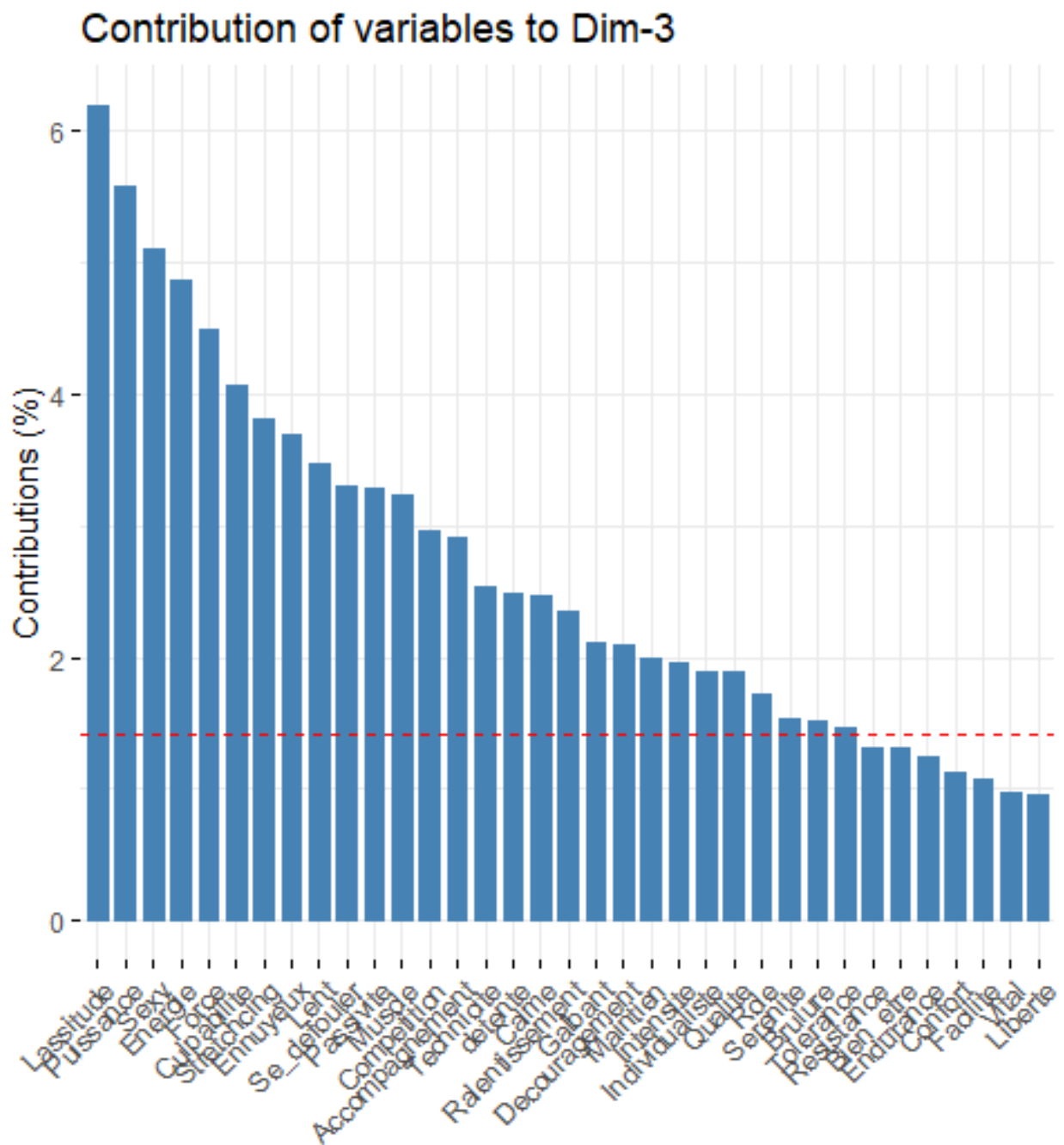


Figure 6:



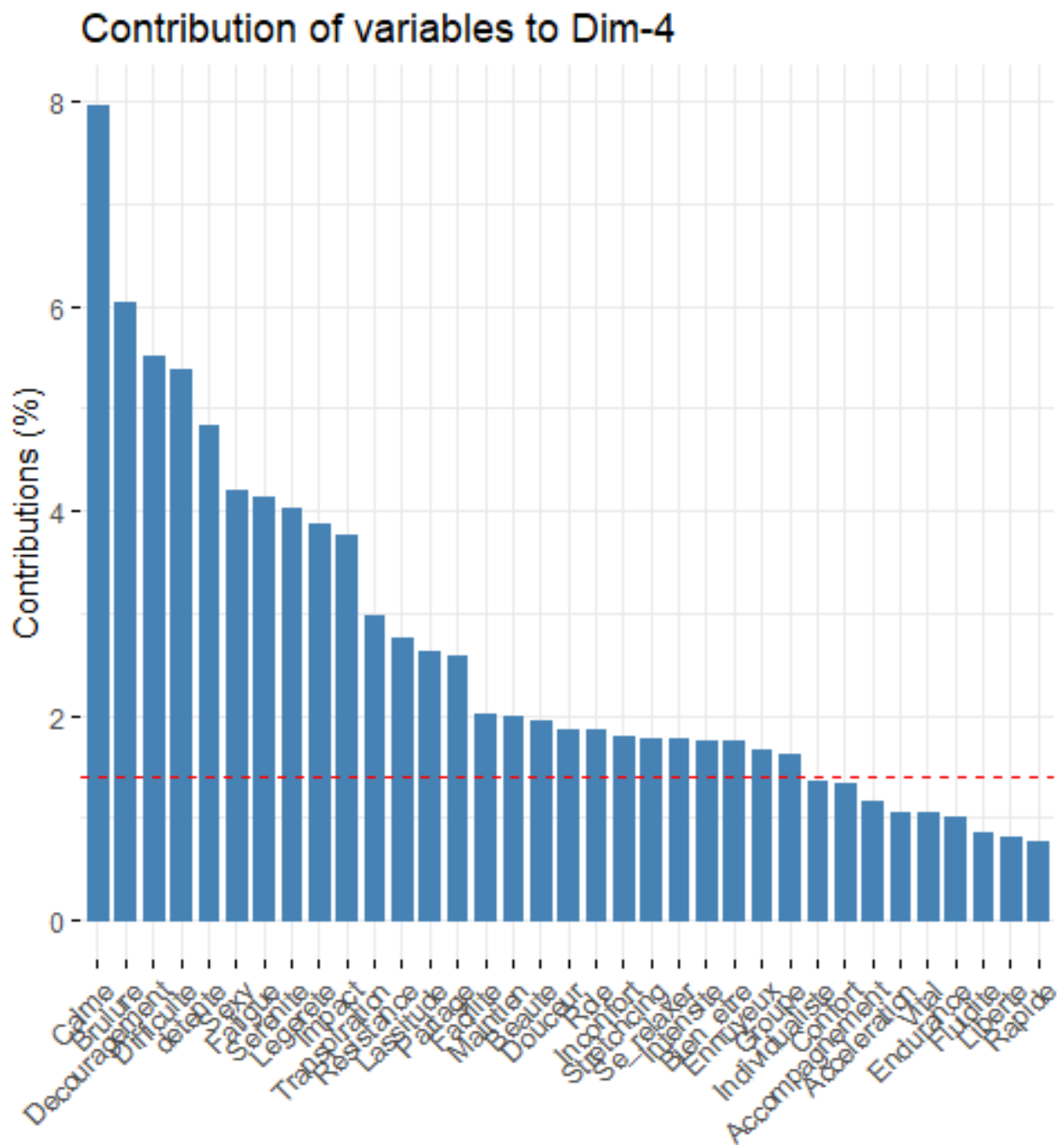
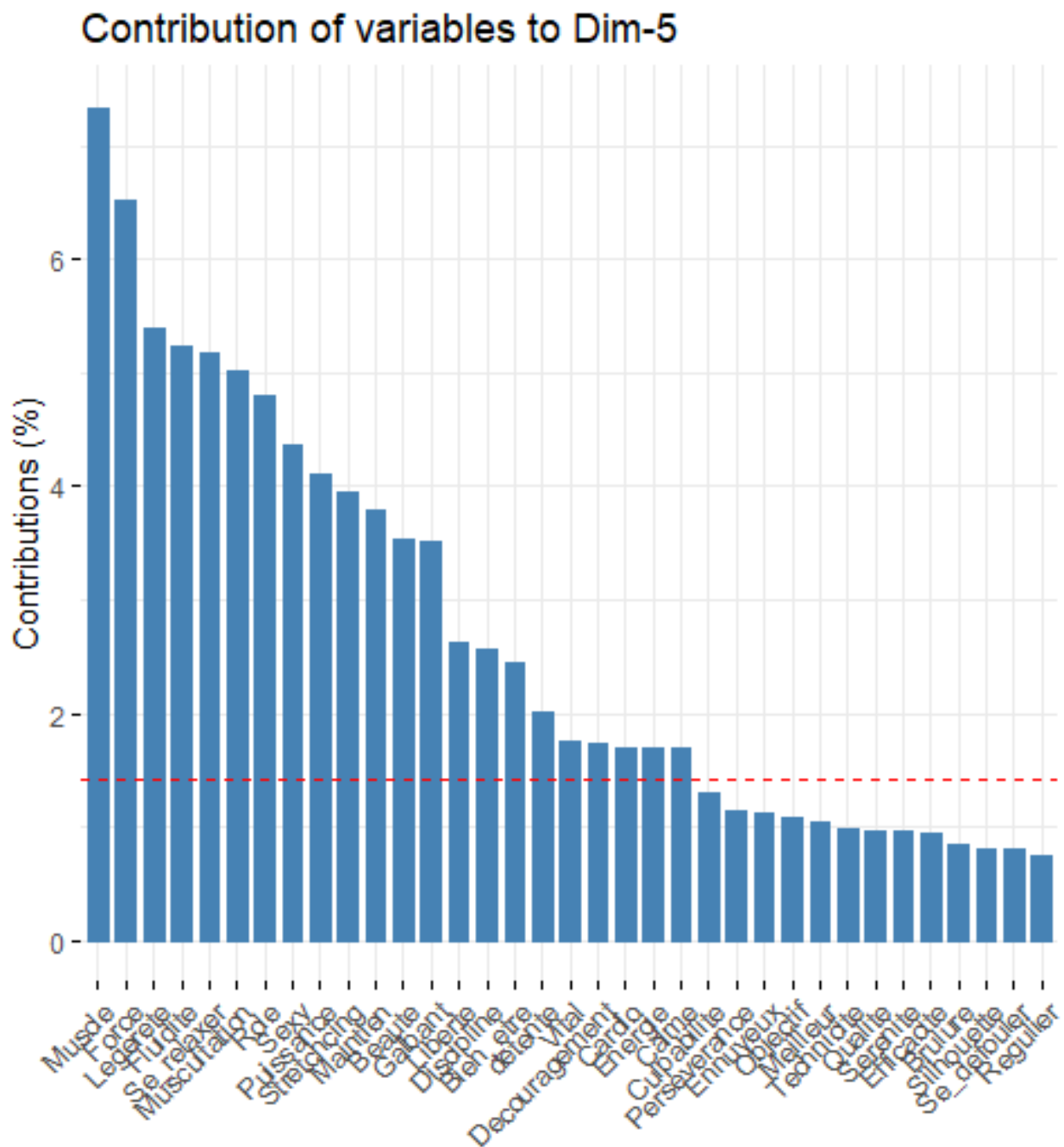


Figure 7:



**Figure 8:**

Les variables les moins importants sont les suivantes : "Recuperation" et "Facilit".

## 4.2 Hierarchical Clustering on Principal Components

Pour réaliser le clustering, nous allons utiliser Hierarchical Clustering on Principal Components (HCPC). Cette méthode permet de combiner les trois méthodes standards utilisées dans les analyses de données multivariées :

- Méthodes en composantes principales (PCA, CA, MCA, FAMD, MFA),
- Regroupement hiérarchique et
- Clustering de partitionnement, en particulier la méthode des k-moyennes.

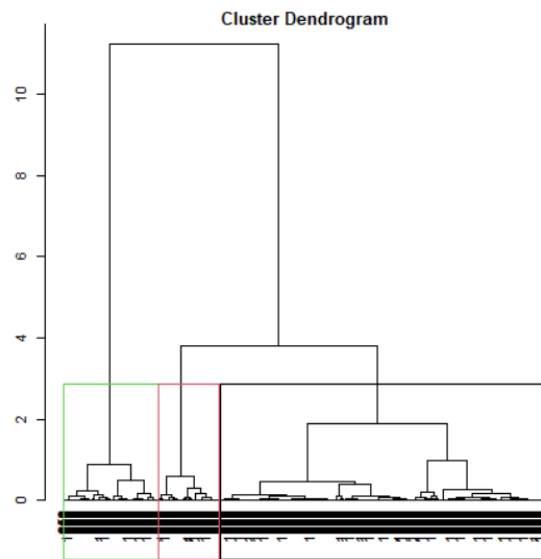
L'algorithme de la méthode HCPC a 4 principales étapes :

- 1) Effectue une ACP. Choisir le nombre de dimensions à retenir en spécifiant l'argument `ncp`. Dans notre cas, la valeur est 5.
- 2) Applique la classification hiérarchique sur le résultat de l'étape 1.
- 3) Choisir le nombre de groupes en fonction du dendrogramme obtenu à l'étape 2. Un partitionnement initial est effectué. Dans notre cas, le nombre de groupes est 3.
- 4) Effectue le k-means pour améliorer le partitionnement initial obtenu à l'étape 3.

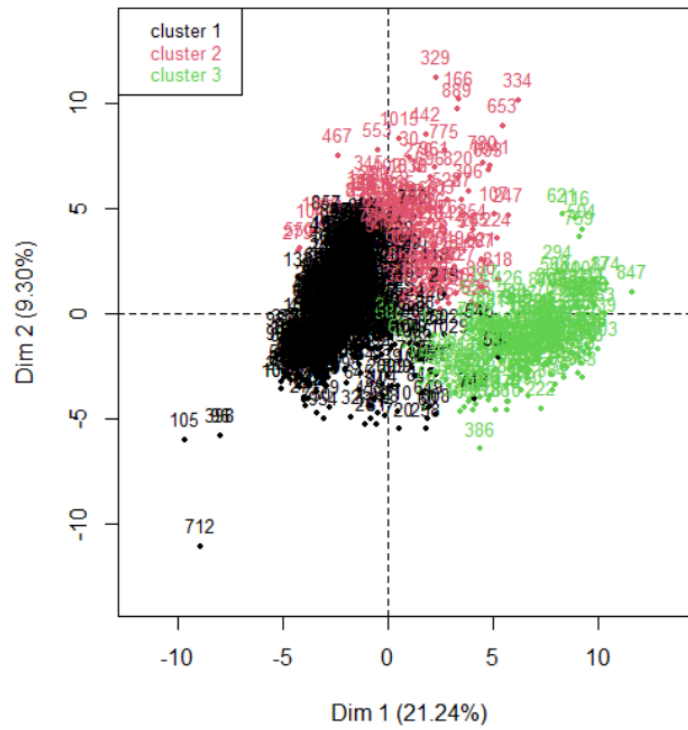
Voici les lignes de codes principales pour :

```
res.pca <- PCA(data_base , ncp = 5 , graph = TRUE)
res.hcpc <- HCPC(res.pca , nb.clust=3, consol=FALSE, graph=TRUE)

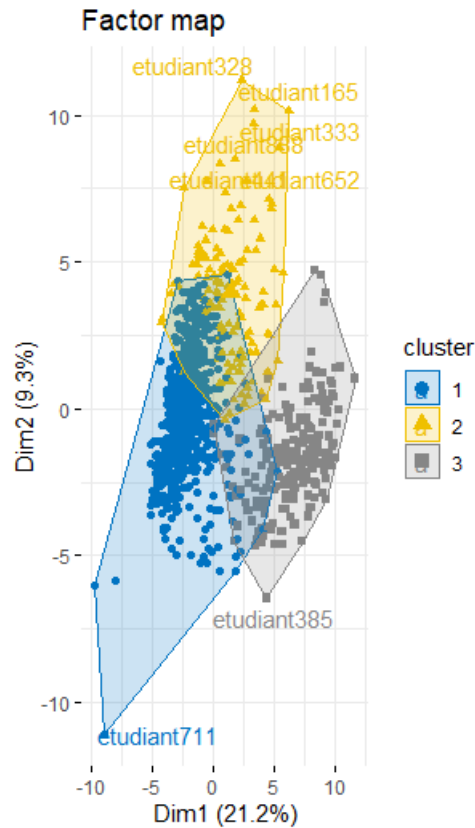
plot(res.hcpc , choice = "tree")
plot(res.hcpc , choice = "map", draw.tree = FALSE)
plot(res.hcpc , choice = "3D.map")
catdes(res.hcpc$data.clust , ncol(res.hcpc$data.clust))
```



**Figure 9:** Hierarchical tree



**Figure 10:** Ascending Hierarchical Classification of the individuals



**Figure 11:** Ascending Hierarchical Classification of the individuals

### 4.3 Results of HCHC

Le cluster 1 est constitué d'individus partageant :

- des valeurs élevées pour les variables Galbant, Culpabilité, Ennuyeux, Stretching, Sécurité, Découragement, Ralentissement, Lassitude, Inconfort et Passivité (les variables sont triées par ordre décroissant).

- des valeurs faibles pour des variables comme Progression, Transpiration, Performance, Actif, Challenge, Plaisir, Objectif, Persévérance, Confiance en soi et Cardio (les variables sont triées par ordre croissant).

Le cluster 2 est constitué d'individus partageant :

- des valeurs élevées pour des variables comme Se défouler, Puissance, Compétition, Technique, Qualité, Énergie, Confort, Muscle, Force et Intensité (les variables sont triées par ordre décroissant).

- des valeurs faibles pour les variables Sexy, Meilleur, Calme et Vital (les variables sont triées par ordre croissant).

Le cluster 3 est constitué d'individus partageant :

- des valeurs élevées pour des variables comme Progression, Actif, Performance, Challenge, Cardio, Partage, Plaisir, Dépassement, Rapide et Efficacité (les variables sont triées par ordre décroissant).

- des valeurs faibles pour des variables telles que Confort, Sécurité, Galbant, Douceur, Ennuyeux, Force, Maintien, Qualité, Beauté et Inconfort (les variables sont triées par ordre croissant).

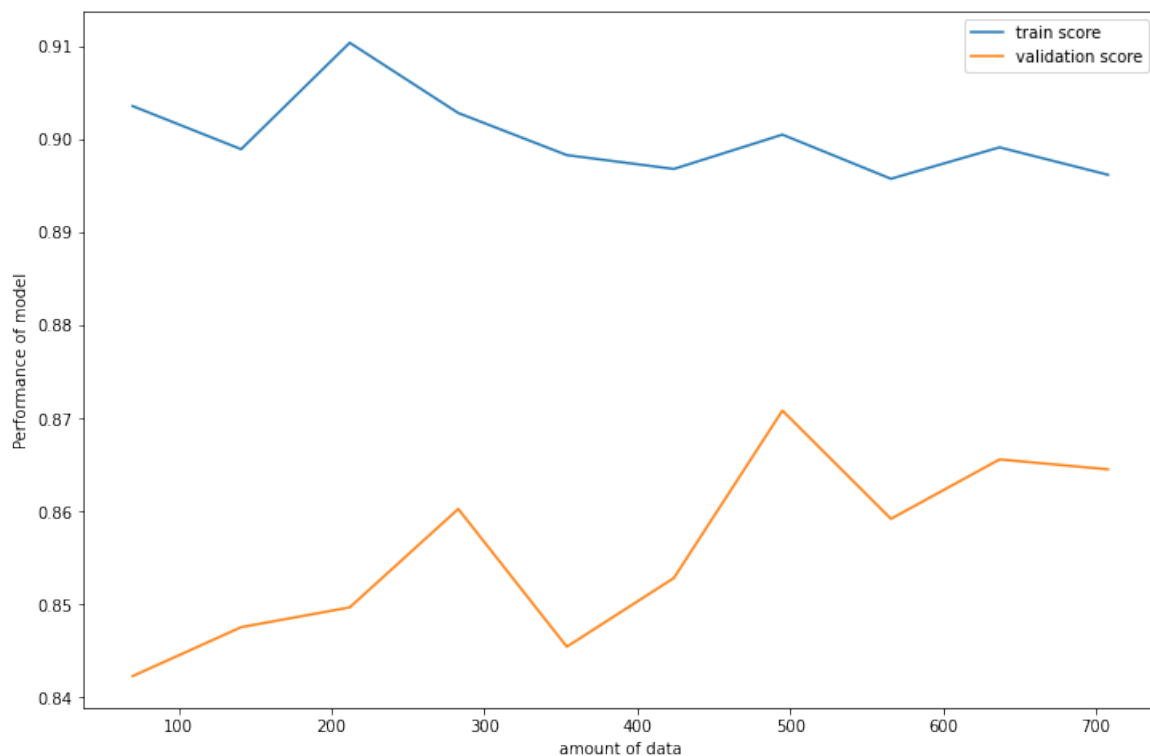
## 5 Classification

Nous avons sparé nos données en 3 parties : - 80 % des données pour le choix de l'algorithme de sélection . - 20 % des données pour tester le modèle final et éventuellement sélectionner les colonnes les plus importantes.

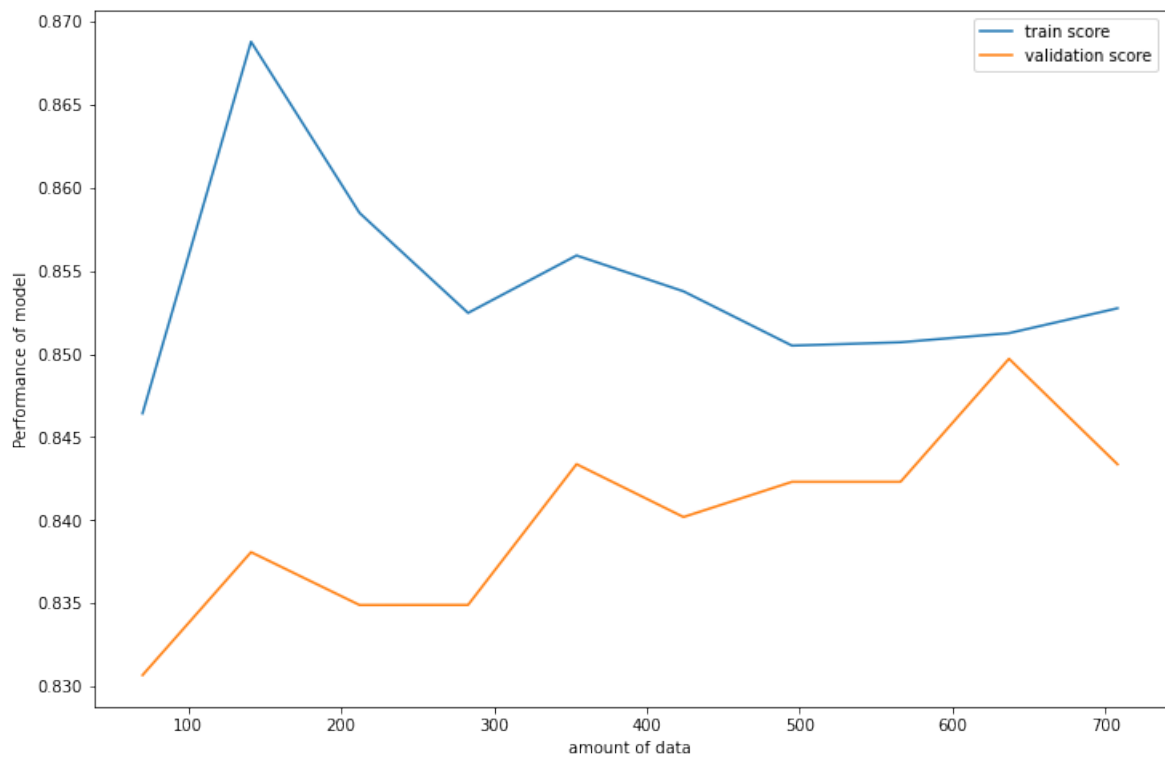
### 5.1 Choice of Classification algorithm

4 classificateurs multi-classes sont utilisés : le classificateur vecteurs de soutien (SVC), classificateur vecteur de support linéaire (LSVC), k-nearest neighbors (KNN) et régression logistique (logreg). (KNN) et régression logistique (logreg).

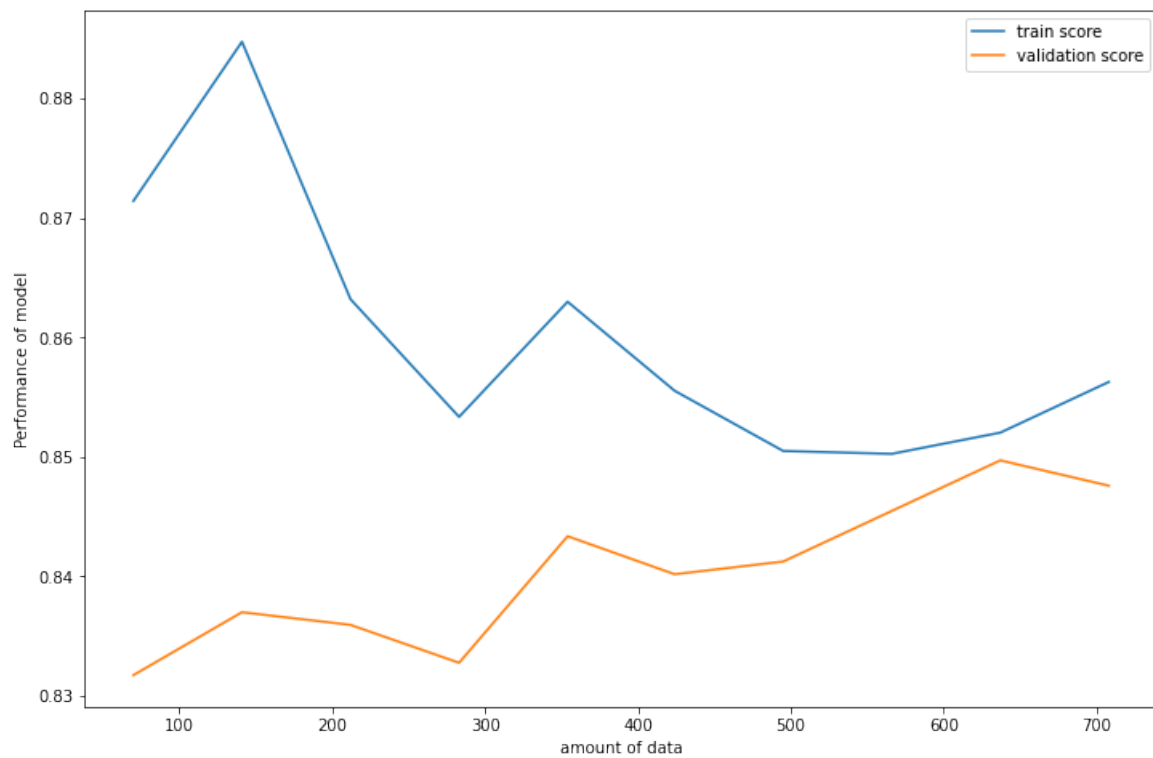
voici le graphique de la performance de chaque modèle :



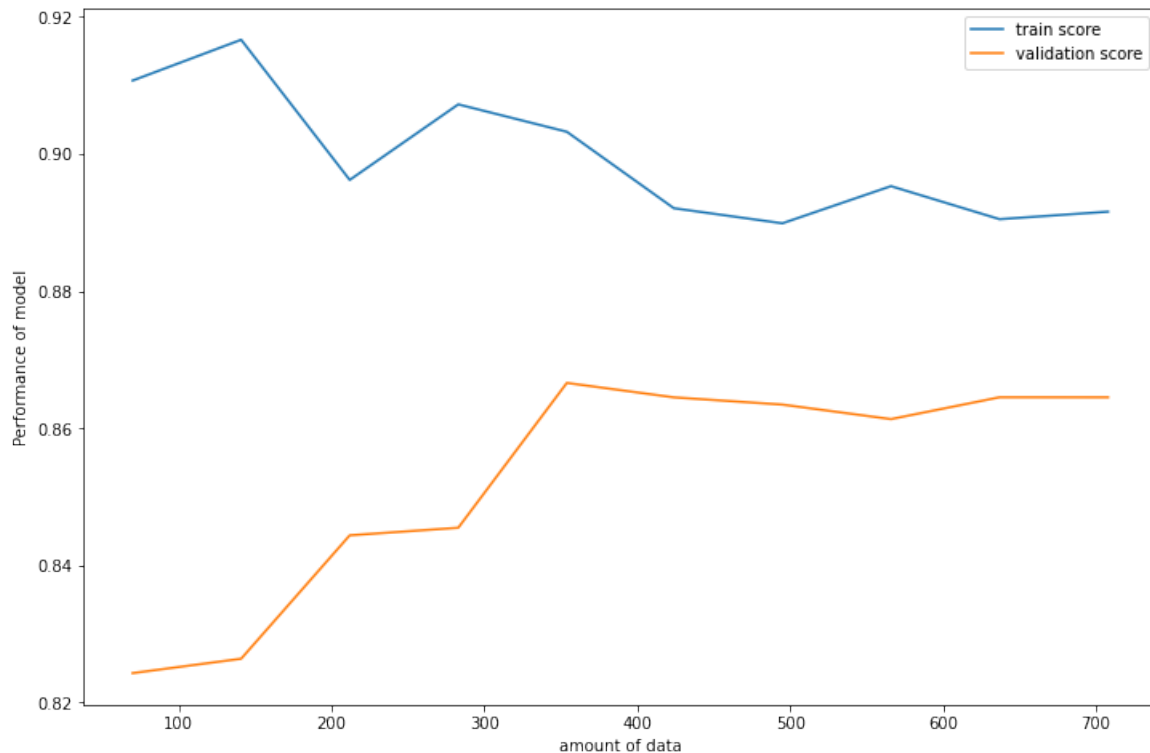
**Figure 12:** KNN learning curve



**Figure 13:** logreg learning curve



**Figure 14:** LSVC learning curve



**Figure 15:** SVC learning curve

D'après les graphiques ci-dessous, la performance du modèle de Logistic Regression et SVC soit plus stable et meilleur que les autres modèles. On peut dire que les deux modèles ne sont pas en overfitting (score train et score val sont proches) contrairement aux 2 autres.

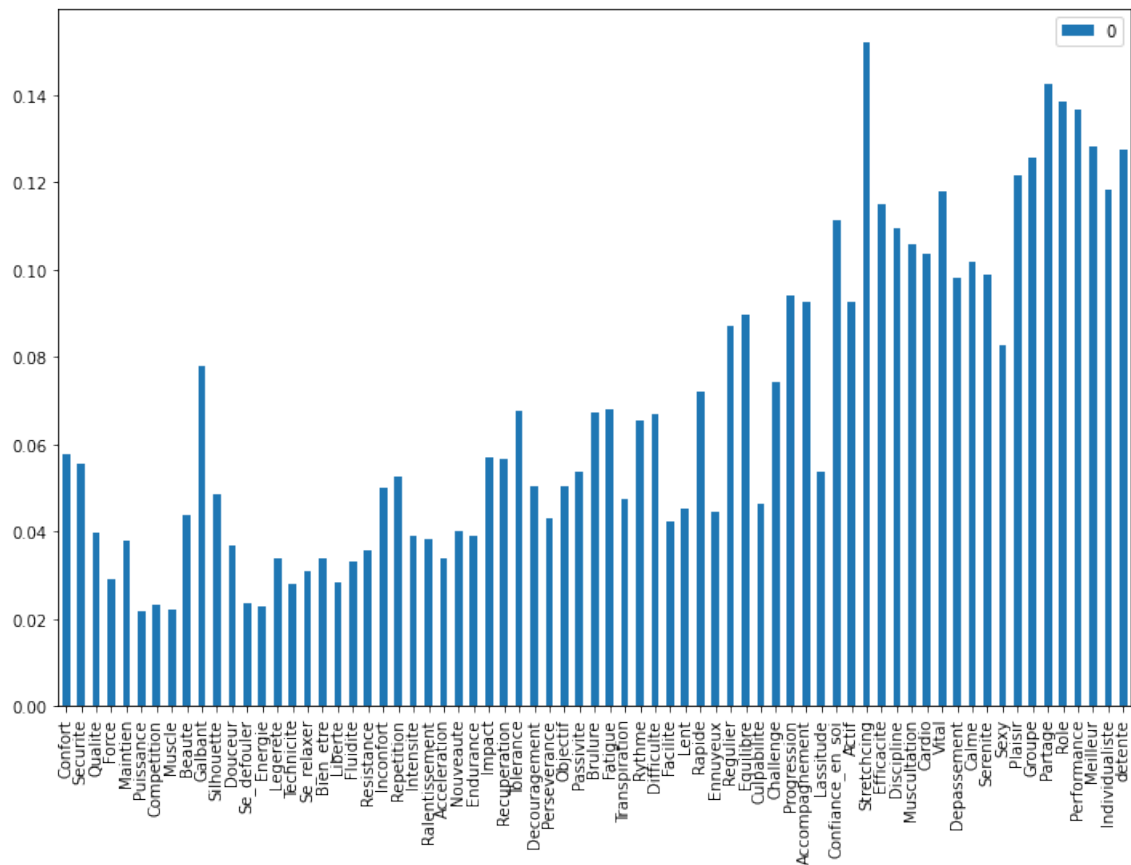
Nous avons utilisé GridSearchCV pour optimiser les hyperparamètres du modèle logistic Regression. GridSearchCV nous permet de les meilleurs hyperparamètres en comparant les différentes performances de chaque combinaison grâce à la technique de cross-validation.

La cross-validation consiste à découper le jeu de données en  $k$  parties égales (ici  $k = 5$ ). À tour de rôle, chacune des  $k$  parties est utilisée comme jeu de test. Le reste (autrement dit, l'union des  $k-1$  autres parties) est utilisé pour l'entraînement.

## 5.2 Feature selection

Le graphique ci-dessous nous montre la variance de chaque feature. 4 candidats le seuil se démarquent : 0.8, 0.06, 0.04 et 0.02.





**Figure 16:** variance of each features

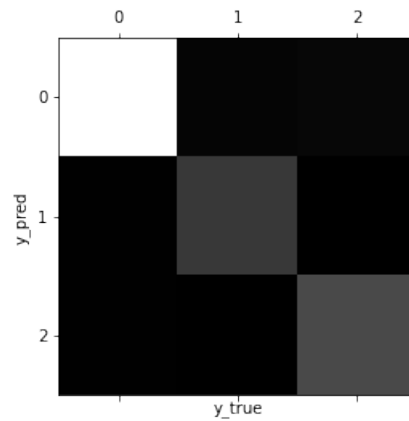
Pour limiter les valeurs inférieures à ce seuil, la fonction `VarianceThreshold` de `scikit-learn` est utilisée.

Au final, le seuil fixé à 0.02 donne de meilleurs résultats. Aucune colonne n'a été supprimée.

### 5.3 Results final model

	predicted: classe 0	predicted: classe 1	predicted: classe 2
Actual: classe 0	133	3	4
Actual: classe 1	1	30	0
Actual: classe 2	1	0	38

**Figure 17:** confusion matrix



**Figure 18:** confusion matrix

Nous obtenons des résultats satisfaisants : seulement 9 valeurs mal placées, la valeur de 0.964 et 0.94768 pour le recall\_score et le f1score.

## 6 Conclusion

L'objectif principal du projet était de faire du clustering sur nos données. Grâce au clustering HCPC, nous pouvons distinguer 3 types d'étudiants :

- les dmotifs, ceux qui recherchent le bien-être et la simplicité dont certaines variables caractéristiques du cluster sont : Culpabilité, Ennuyeux, Découragement, Ralentissement, Inconfort
- Dans le deuxième groupe, on a ceux qui aiment les sports de combat comme la lutte, la boxe ou le MMA. Il est caractérisé par ses mots : Puissance, Compétition, Technique, Qualité, Énergie, Muscle, Force et Intensité
- Le dernier groupe se distingue par ses mots : Progression, Performance, Challenge, Cardio, Partage, Dépassement, Rapide et Efficacité. On retrouve ceux qui apprécient la course à pied et les activités de nature.

En ce qui concerne la classification, l'un des meilleurs algorithmes (SVM) a donné la valeur 0.9549 de précision, 0.9 de rappel et 0.9255 de f1-score. (Par contre, l'ensemble de données de ce projet n'est pas assez grand.)

## References

- [1] <https://www.cairn.info/revue-staps-2018-2-page-99.htm>
- [2] <https://solidarites-sante.gouv.fr/prevention-en-sante/preserver-sa-sante/article/activite-physique-et-sante>
- [3] <https://e3s.unistra.fr/equipe/presentation/>
- [4] <https://scikit-learn.org>
- [5] [https://fr.wikipedia.org/wiki/Analyse\\_en\\_composantes\\_principales](https://fr.wikipedia.org/wiki/Analyse_en_composantes_principales)

- [6] <http://www.sthda.com/english/>
- [7] <http://www.sthda.com/english/articles/31-principal-component-methods-in-r-practical-117-hcpc-hierarchical-clustering-on-principal-components-essentials/#algorithm-of-the-hcpc-method>
- [8] <https://husson.github.io/teaching.html>
- [9] <https://www.youtube.com/c/MachineLearnia>
- [10]