

Rendu TP1 : Apprentissage non supervisé

Analyse en composantes principales

On centre les données pour garder leurs variations et non leurs valeurs absolues.

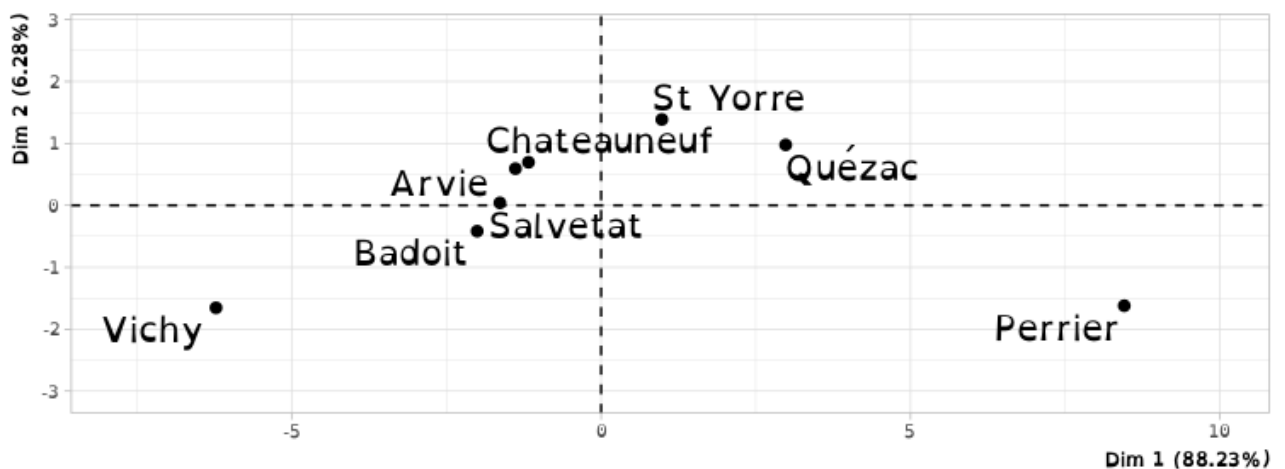
Sur la matrice de corrélation, on voit que les saveurs **alcaline** et **sucrée** indiquent une caractéristique commune et que l'**acide** et l'**amère** indiquent une caractéristique proche. À l'inverse la relation entre les saveurs **amère** et **sucrée**, **amère** et **alcaline**, **acide** et **alcaline** indique des caractéristiques très distinctes.

La matrice des distances entre les individus nous indique les individus qui ont les données les moins similaires. Par exemple, **Vichy** et **Perrier** sont très différentes avec une valeur de 14.7. On peut le vérifier directement dans les données brutes.

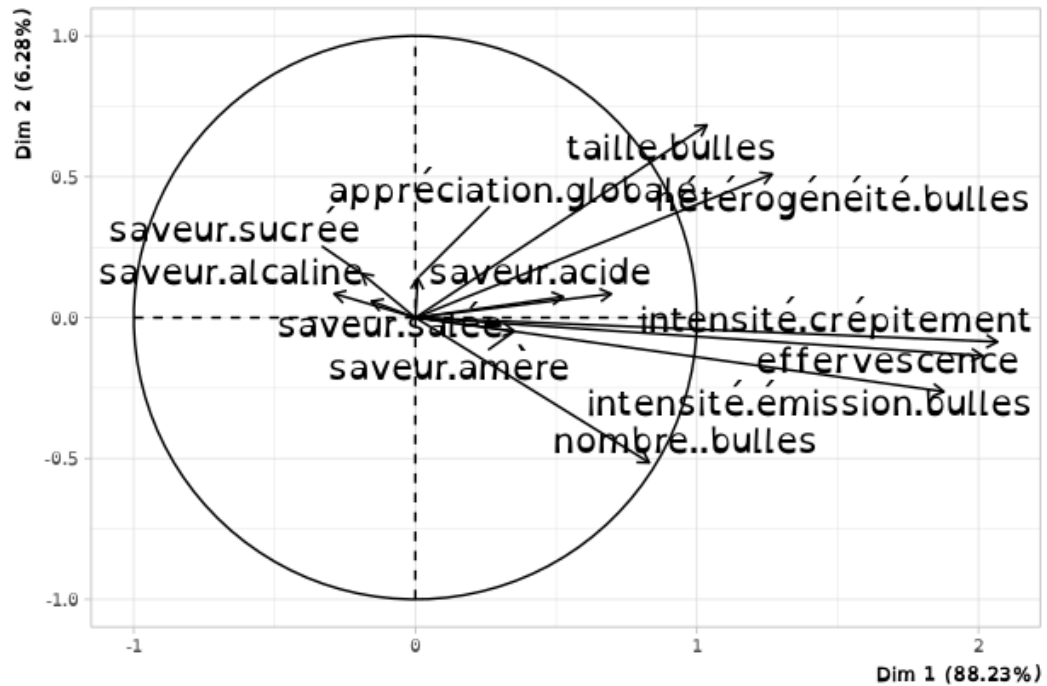
	saveur.amère	saveur.sucrée	saveur.acide	saveur.salée	saveur.alcaline	appréciation.globale	intensité.émission.bulles
St Yorre	-0.5072997	1.03143052	-0.2709005	1.6833150	0.60915786	-0.5651411	0.012321
Badoit	0.3043798	-0.55538567	-0.6043166	-0.6059934	-0.20305262	-0.5651411	-0.381980
Vichy	-1.5218990	0.39670405	-1.6045647	1.1446542	1.15063151	-0.9129202	-1.219873
Quézac	0.5072997	-0.55538567	1.2294716	-0.6059934	-0.74452627	1.5215337	0.603775
Arvie	-1.1160592	1.34879376	-0.1041925	0.0673326	1.15063151	-0.5651411	-0.381980
Chateauneuf	0.1014599	0.07934081	-0.1041925	0.0673326	0.06768421	0.8259754	-0.283405
Salvetat	0.7102195	0.07934081	-0.1041925	-1.4139846	-0.20305262	1.1737546	-0.480556
Perrier	1.5218990	-1.82483862	1.5628877	-0.3366630	-1.82747358	-0.9129202	2.131698

Pour l'ACP

- Sur la figure des individus, on voit des individus groupés au centre et quelques individus éloignés.

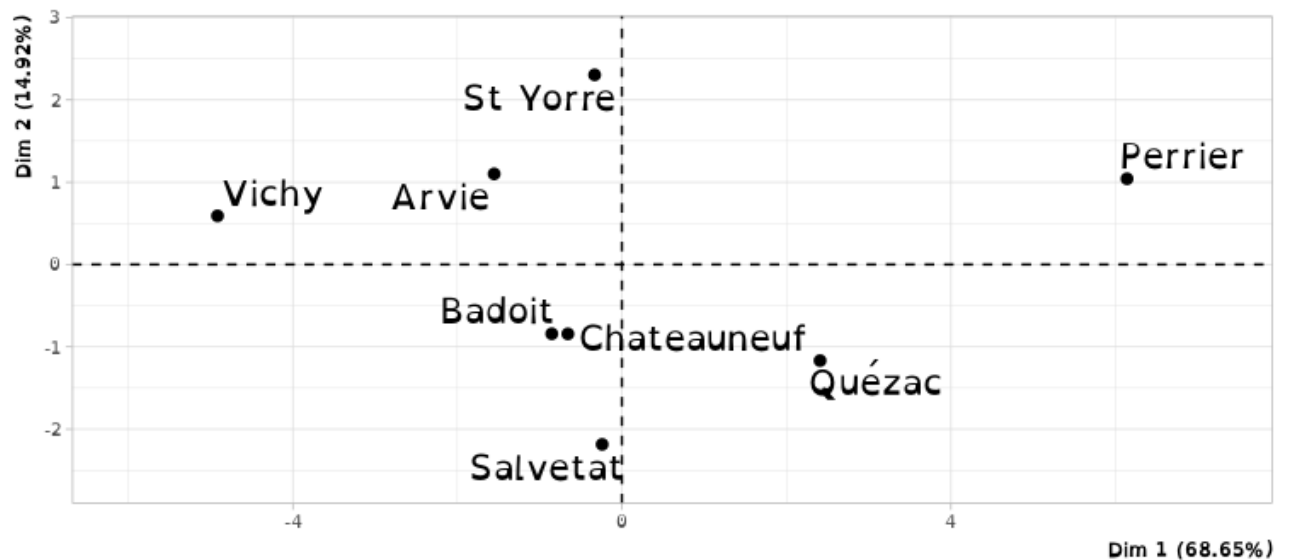


- Sur la figure des variables, on voit que des vecteurs sortent du cercle °°.

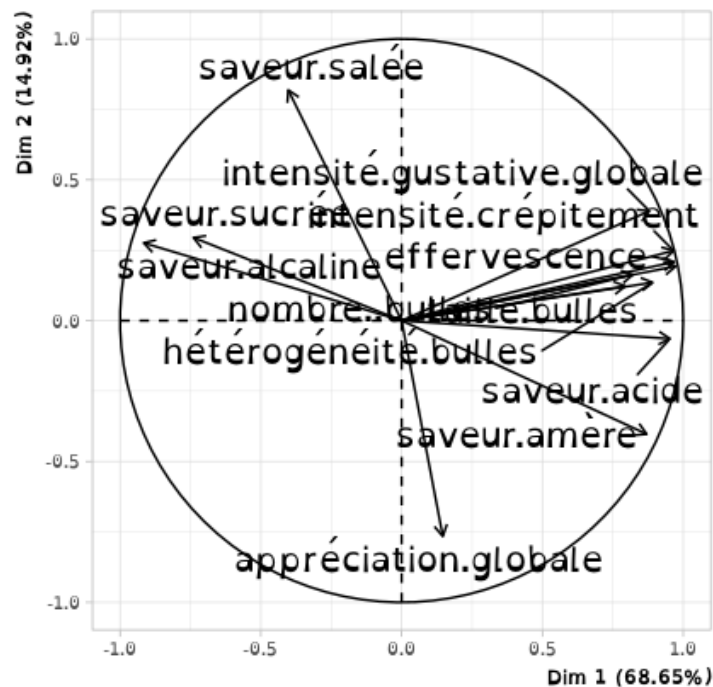


Pour l'ACP normé

- Sur la figure des individus, on voit que les individus sont plus dispersés, mais quelques individus restent éloignés comme **Vichy** et **Perrier**.



- Sur la figure des variables, on voit que des vecteurs ne sortent plus du cercle °~°.



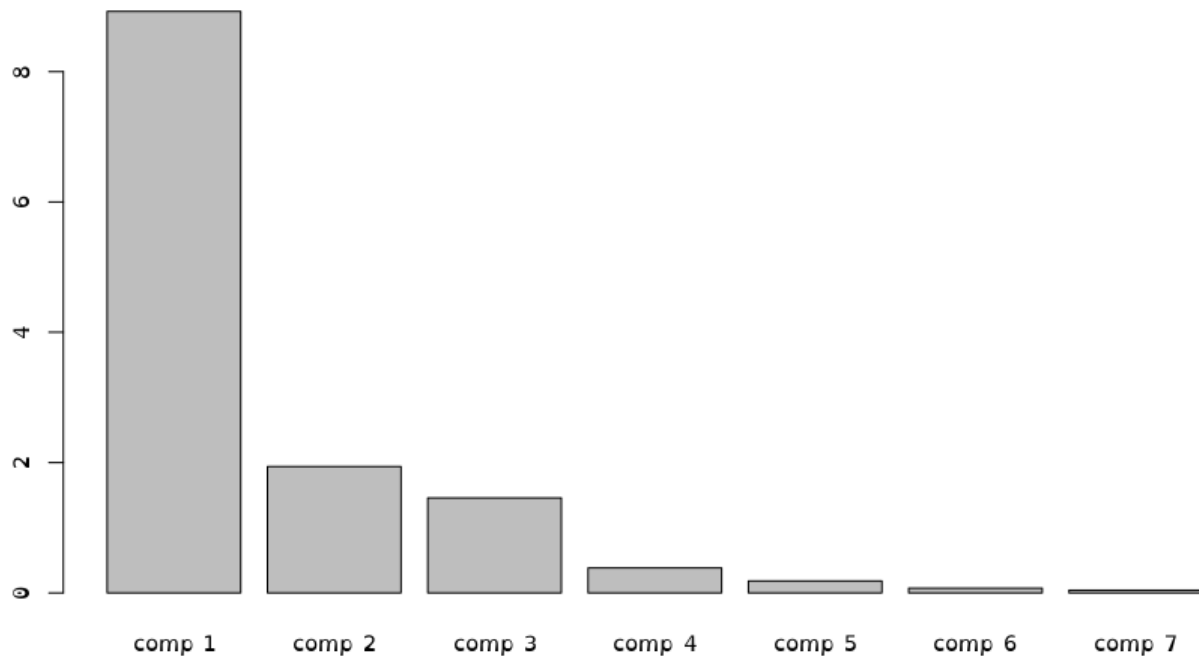
Choix du nombre de composantes à retenir

Avec la **règle de Kaiser** on garde **3** composantes, car les **3** premières valeurs propres sont supérieures à **1**.

Les valeurs propres sont nommées λ_d avec d la dimension (c'est la variance empirique). La valeur propre par dimension :

- Avec **1** composante on couvre **68.65%** des données significatives et une valeur propre de $\lambda_1=8.9$.
- Avec **2** composantes on couvre **83.57%** des données significatives et une valeur propre de $\lambda_2=1.9$.
- Avec **3** composantes on couvre **94.77%** des données significatives et une valeur propre de $\lambda_3=1.5$.
- Avec **4** composantes on couvre **97.73%** des données significatives et une valeur propre de $\lambda_4=0.3$.

Avec la règle du coude on garde **2** composantes. Visuellement, on voit que la cassure est dès la **2**ème composante.



Par le calcul, on tombe sur 2 aussi, on arrête le calcul dès que δ_x est négatif puis on compte le nombre de δ_x .

Calcul des différences premières :

- $\epsilon_1 = (\lambda_1 - \lambda_2)$, $\epsilon_2 = (\lambda_2 - \lambda_3)$, $\epsilon_3 = (\lambda_3 - \lambda_4)$
- $\epsilon_1 = (8.92 - 1.93)$, $\epsilon_2 = (1.93 - 1.45)$, $\epsilon_3 = (1.45 - 0.38)$
- $\epsilon_1 = 7.99$, $\epsilon_2 = 0.48$, $\epsilon_3 = 1.07$

Calcul des différences secondes :

- $\delta_1 = (\epsilon_1 - \epsilon_2)$, $\delta_2 = (\epsilon_2 - \epsilon_3)$
- $\delta_1 = (7.99 - 0.48)$, $\delta_2 = (0.48 - 1.07)$
- $\delta_1 = 7.51$, $\delta_2 = -0.59$

On a une bonne qualité de la projection et contribution des individus avec 3 dimensions, car ils contribuent tous à plus de 71%, ce qui est mieux que 2 dimensions où la contribution la plus basse est de 55% :

```
resnorm$ind$cos2[,1]+resnorm$ind$cos2[,2]
St Yorre      Badoit      Vichy      Quézac      Arvie Chateauneuf      Salvetat      Perrier
0.6912732    0.3777300    0.8882563    0.7323157    0.6012887    0.5518255    0.8856256    0.9360553
# Et les 3 premières ?
resnorm$ind$cos2[,1]+resnorm$ind$cos2[,2]+resnorm$ind$cos2[,3]
St Yorre      Badoit      Vichy      Quézac      Arvie Chateauneuf      Salvetat      Perrier
0.9405832    0.7955333    0.9777641    0.9271582    0.7145190    0.7536751    0.8925276    0.9988988
```

On a une bonne qualité de la projection et contribution des variables avec 3 dimensions, car elles contribuent toutes à plus de 84%, ce qui est mieux que 2 dimensions où la contribution la plus basse est de 60% :

saveur.amère	0.9193755	saveur.sucrée	0.6333039	saveur.acide	0.9122874
saveur.salée	0.8321232	saveur.alcaline	0.9153633	appréciation.globale	0.6073619
intensité.émission.bulles	0.9842741	nombre.bulles	0.7013436	taille.bulles	0.6455454
hétérogénéité.bulles	0.8128302	effervescence	0.9801694	intensité.gustative.globale	0.9252784
intensité.crépitement	0.9945345				

Et les 3 ?

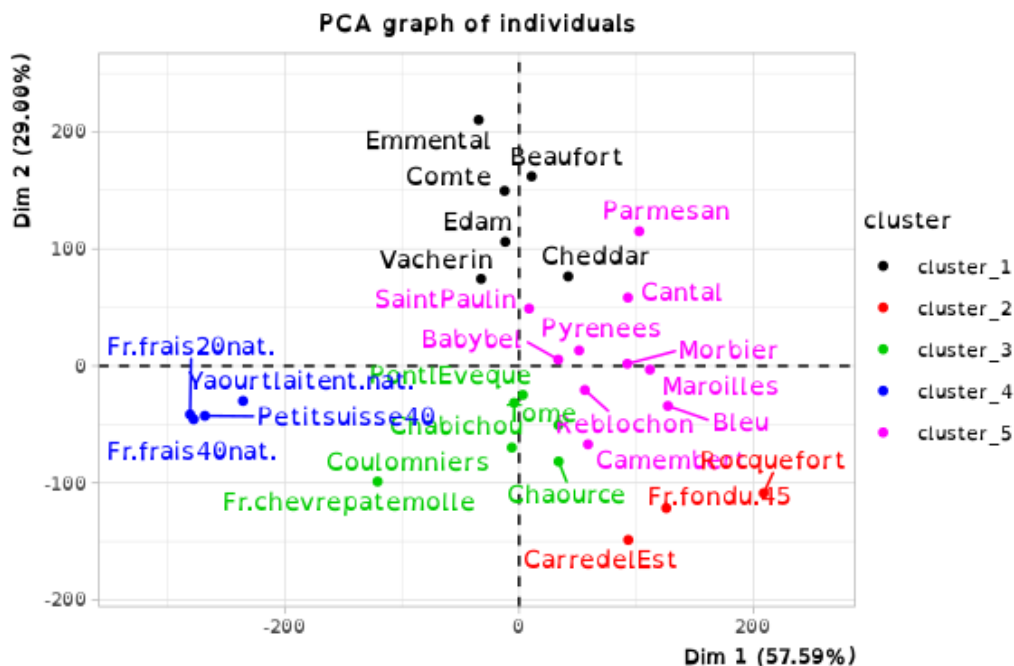
```
resnorm$var$cos2[,1]+resnorm$var$cos2[,2]+resnorm$var$cos2[,3]
```

saveur.amère	0.9362833	saveur.sucrée	0.9247121	saveur.acide	0.9568183
saveur.salée	0.8409480	saveur.alcaline	0.9733153	appréciation.globale	0.9419222
intensité.émission.bulles	0.9917772	nombre.bulles	0.9214411	taille.bulles	0.9139851
hétérogénéité.bulles	0.9842623	effervescence	0.9802432	intensité.gustative.globale	0.9605466
intensité.crépitement	0.9948487				

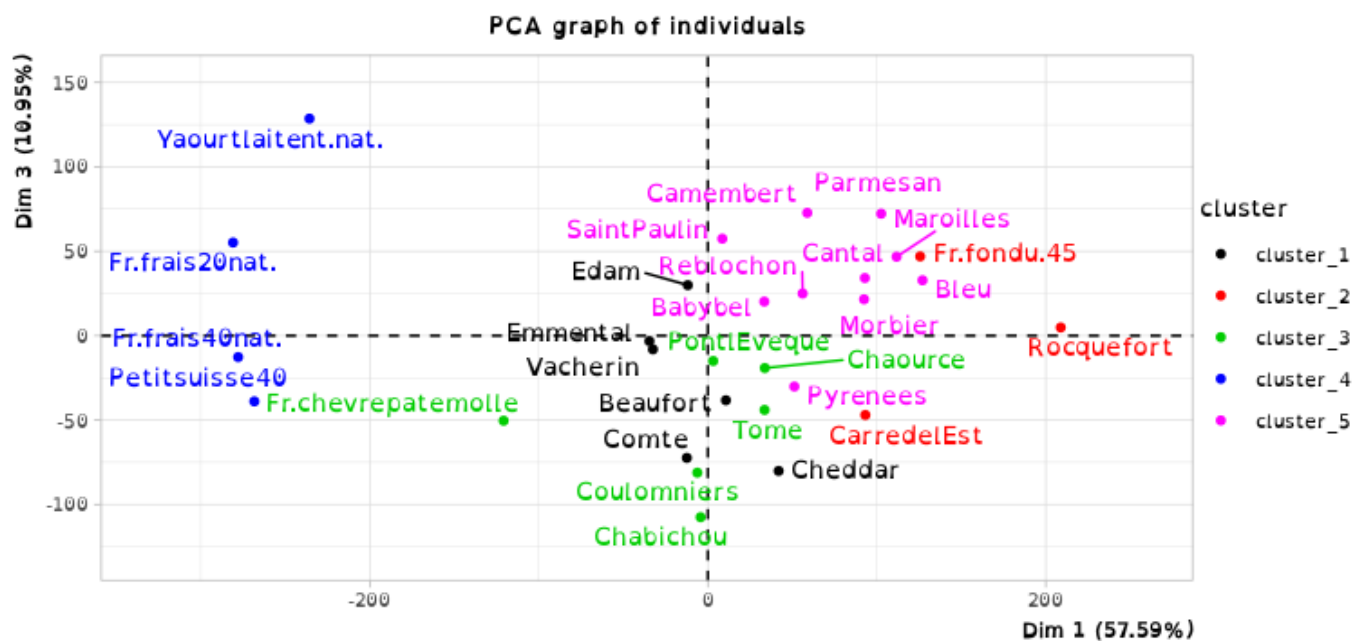
Partitionnement

Les individus sont mieux répartis pour visuellement voir les clusters avec le **k-mean** brute.

K-mean sur les données brute :



K-mean sur les données centrées-réduites :

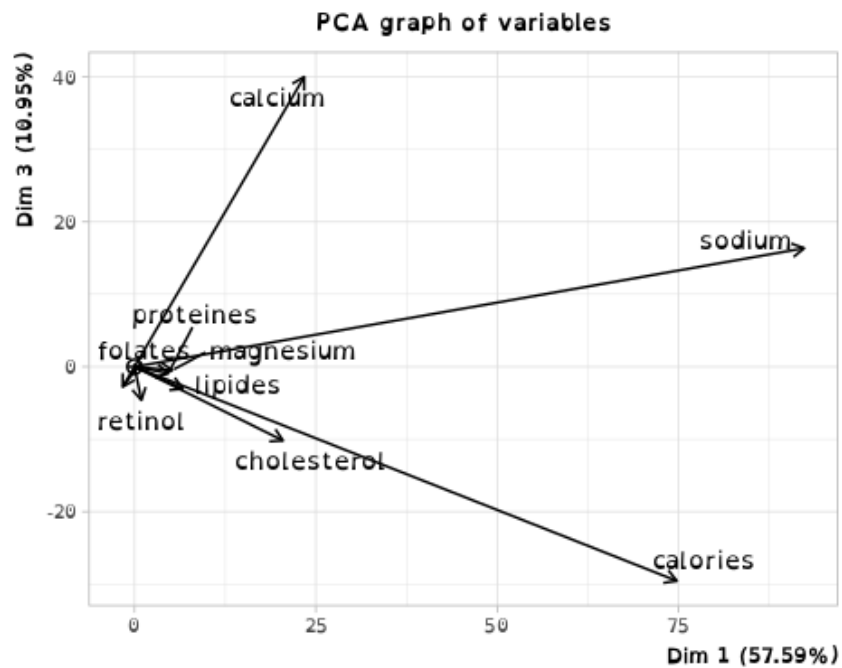


On distingue mieux les variables avec le K-mean centrées-réduites.

K-mean sur les données brute :



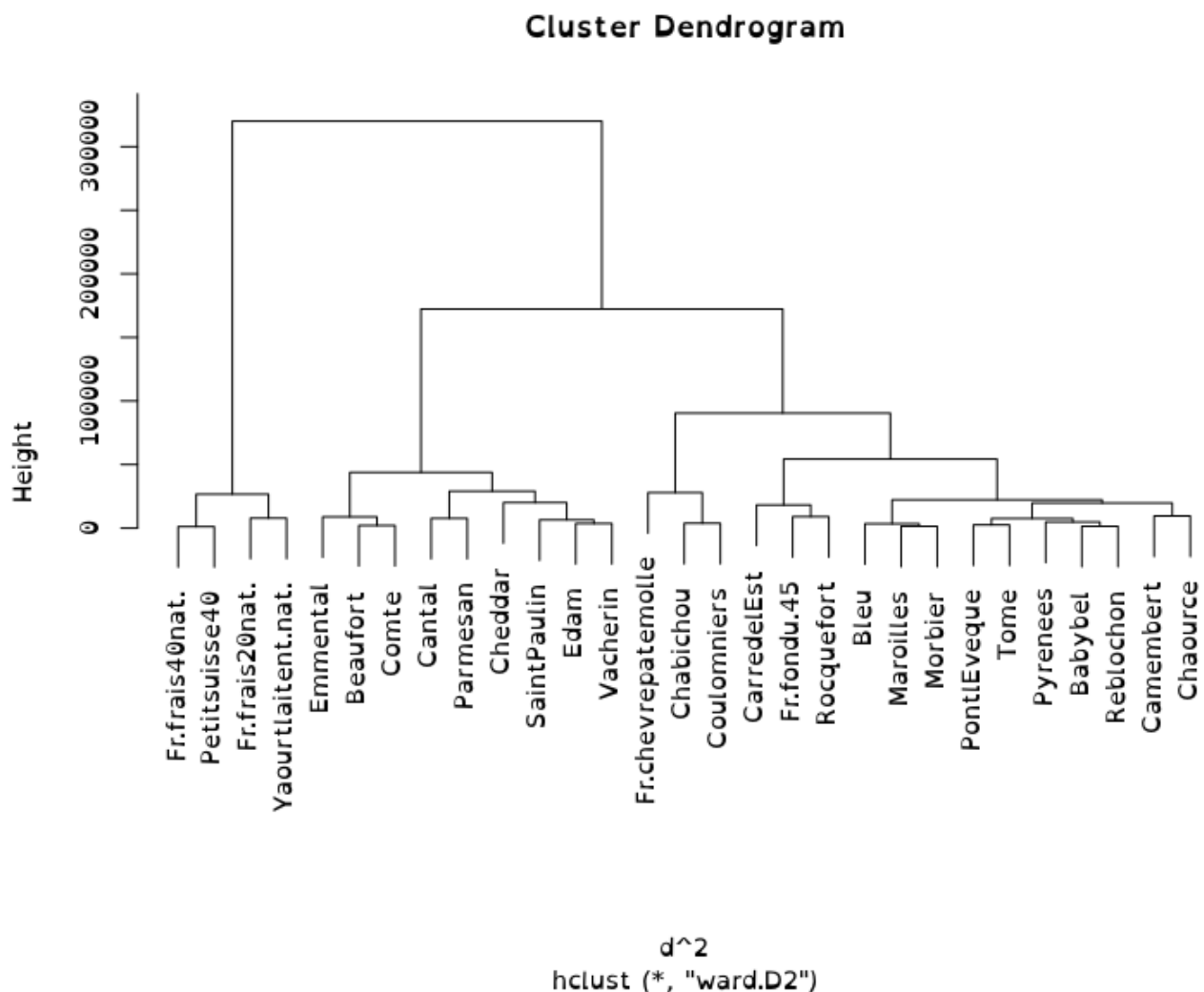
K-mean sur les données centrées-réduites :



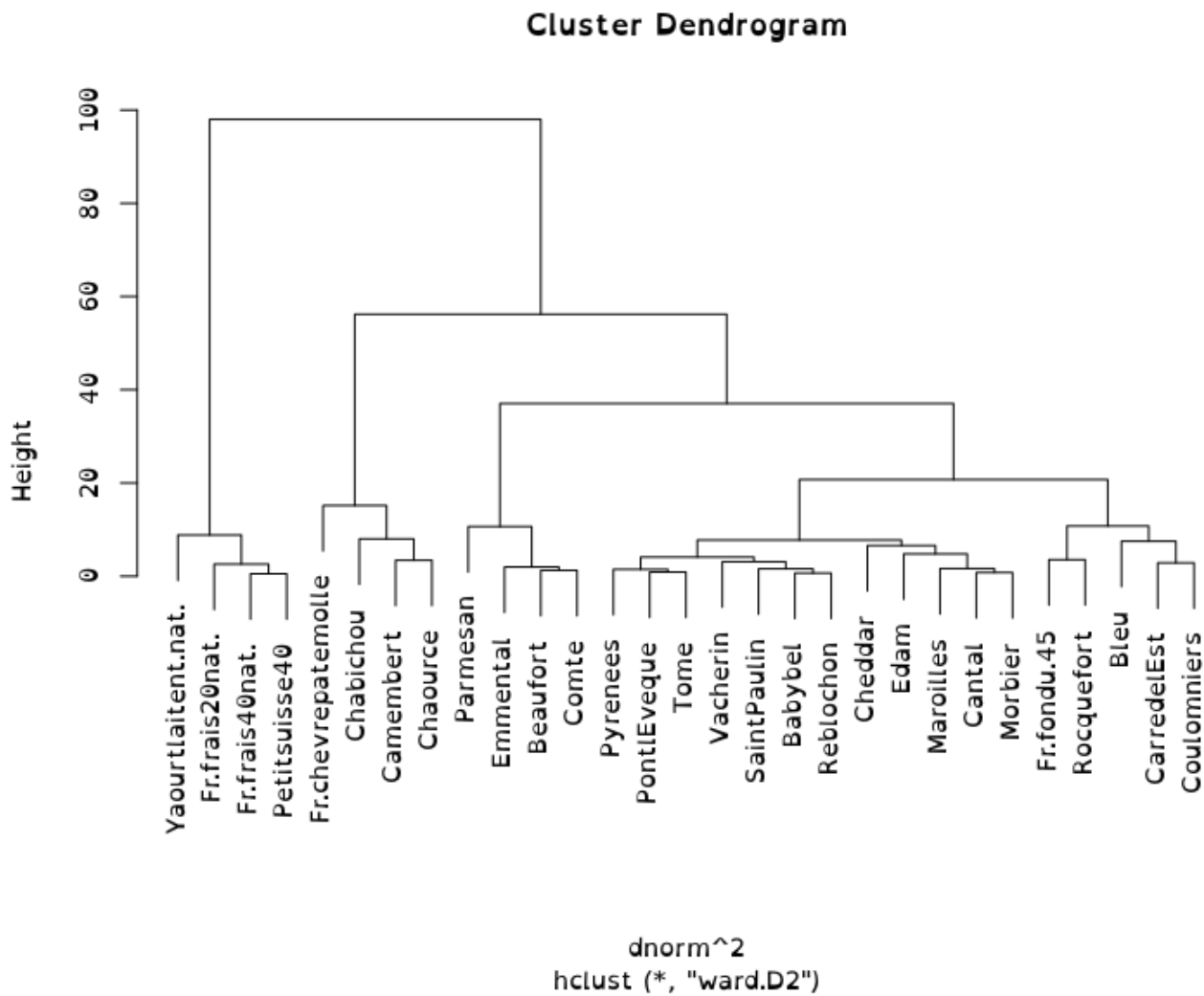
Partitionnement hiérarchique (distance de Ward)

Le découpage est plus distinct avec les données centrées-réduites.

Classification hiérarchique de Ward sur les données brutes :

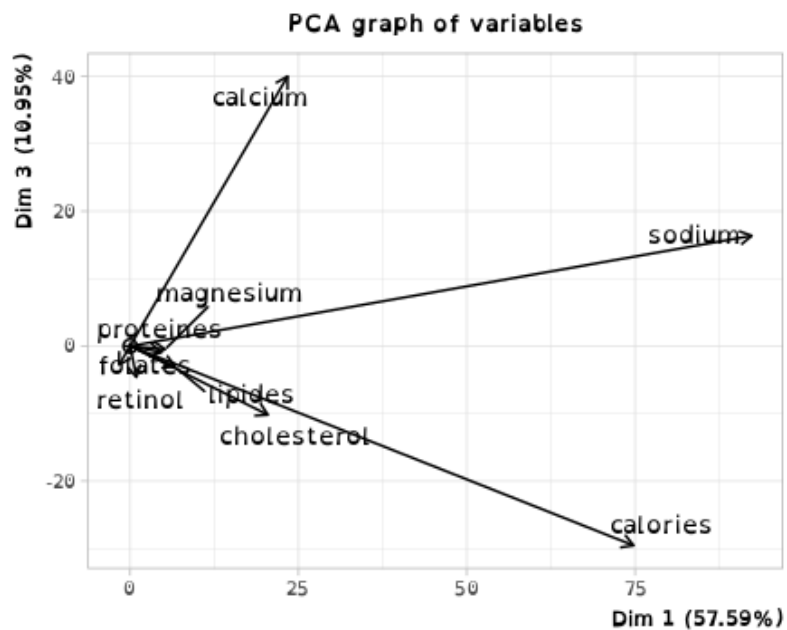
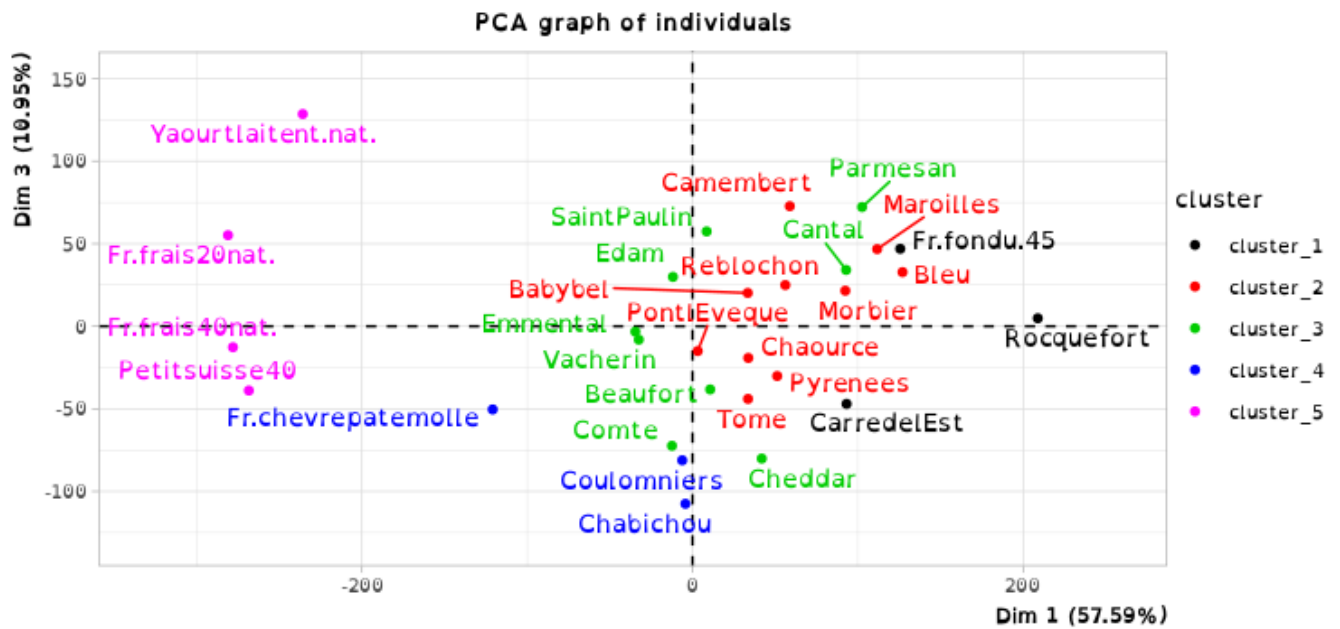


Classification hiérarchique de Ward sur les données centrées-réduites :



Avec la **distance de Ward** sur les données centrées-réduites, on voit que les clusters sont mieux répartis et définis.

distance de Ward sur les données brute :



distance de Ward sur les données centrées-réduites :

