



TD N°2 : Analyse Factorielle des Correspondances (AFC)

Analyse exploratoire

Auteurs : Manelle Nouar, Eliott Bernardou, Mohamed Abdelkader, Sarra Madad, Alfio Chadoin
Enseignants : Myriam Bertrand & Jordan Gonzalez

25 septembre 2022

1 | Elections presidentielles 2012

1. Importer le jeu de données President12.txt. Dans la suite, le jeu de données importe sous R est appelé **president12**. Déclarer le département comme identificateur des individus (nom des cas dans R).

```

1  president12 = read.csv2("/cloud/project/President12.csv")
2  president12 = as.data.frame(president12)
3
4  View(president12)
5  str(president12)
6  dim(president12)
7  names(president12)
8  attach(president12)
9  summary(president12)
10 dim(president12)
11 colnames(president12)
12 row.names(president12)
13
14 rownames(president12) = president12[, 2]
```

Dans ce TP, nous étudions un jeu de données qui représente les présidentielles de 2012. Nous avons 108 observations et 13 variables.

```

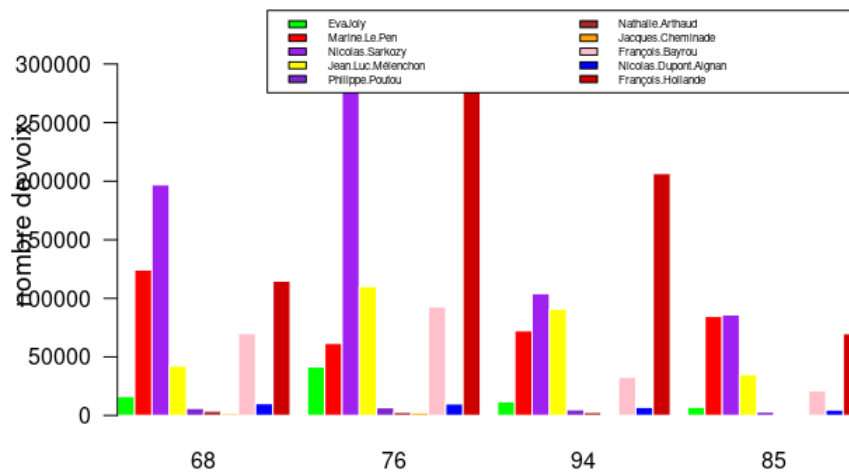
> rownames(president12)
[1] "Ain" "Aisne"
[3] "Allier" "Alpes de Haute.Provence"
[5] "Hautes.Alpes" "Alpes.Maritimes"
[7] "Ardèche" "Ardennes"
[9] "Ariège" "Aube"
[11] "Aude" "Aveyron"
[13] "Bouches.du.Rhône" "Calvados"
[15] "Cantal" "Charente"
[17] "Charente.Maritime" "Cher"
[19] "Corrèze" "Corse.du.Sud"
[21] "Haute.Corse" "Côte.d.Or"
[23] "Côtes d.Armor" "Creuse"
[25] "Dordogne" "Doubs"
```

2. Représenter par un diagramme en batons les nombres de voix obtenues par les différents candidats dans les départements du Bas-Rhin, de Paris, de la Seine Saint- Denis, du Vaucluse.

```

1  color = c("green", "red", "purple", "yellow",
2            "purple3", "brown", "orange", "pink", "blue", "red3")
3
4  barplot(t(president12[c(68, 76, 94, 85), 3:12]),
5          col = color, border="white", horiz=F, beside=T,
6          ylab="nombre de voix", las=1, cex.lab=1.2)
7
8  legend("topright", cex = 0.5, ncol = 2,
9        legend = colnames(president12)[3:12],
10         fill = color)
11
```

Notre histogramme affiche le nombre de voix dans le Bas-Rhin, à Paris, en Seine Saint-Denis et dans le Vaucluse.



3. Construire :

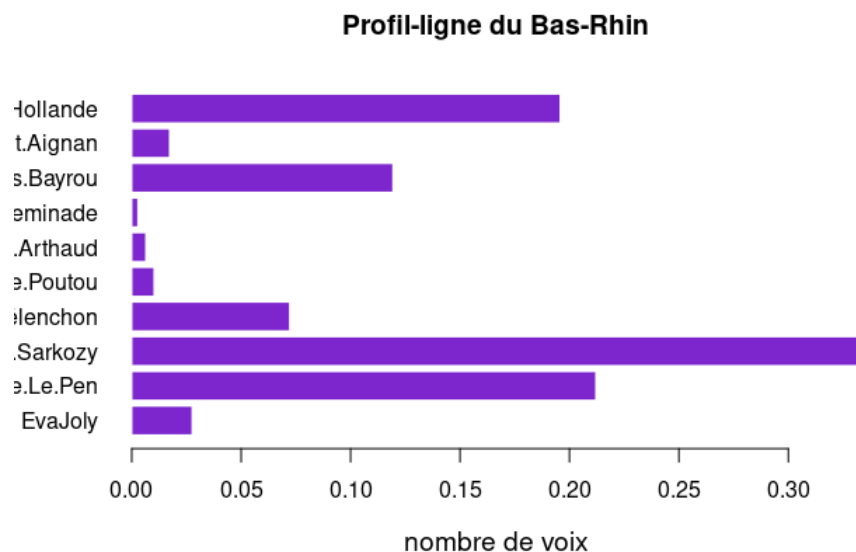
- le profil-ligne du Bas-Rhin (tableau des fréquences).
- le graphique représentant le profil-ligne du Bas-Rhin.

```

1 tab_freq = president12[68, 3:12] / president12[68, 13] # Tableau des
  fréquences
2
3 barplot(as.matrix(tab_freq), names.arg=colnames(president12)[3:12],
4         col="purple3", border="white", main="Profil-ligne du Bas-Rhin",
5         horiz=T, las=1, xlab="nombre de voix", cex.lab=1.2)

```

Nous avons affiché le profil-ligne du Bas-Rhin :



4. **Realiser le test du χ^2 permettant d'étudier le lien de dépendance entre les 96 départements de France métropolitaine et les dix candidats.**
- Extraire des résultats de ce test : la p-value puis les données attendues en cas d'indépendance .
 - Quel est le profil ligne attendu d'un département quelconque sous l'hypothèse d'indépendance ?

```
1 sum(president12[1:96,3:12])
2 min(chisq.test(president12[1:96,3:12])$expected)
3
4 resu.chi2 = chisq.test(president12[1:96,3:12])
5 resu.chi2
```

Ce test permet de vérifier l'absence de lien statistique entre deux variables X et Y. Les deux sont dites indépendantes lorsqu'il n'existe aucun lien statistique entre elles.

Pour effectuer ce test il faut remplir certaines conditions comme, l'effectif total > 50, l'effectif théorique minimum supérieur à 5 et avoir des variables qualitatives.

L'hypothèse nulle (H_0) de ce test est la suivante : les deux variables X et Y sont indépendantes.

En termes de valeur p, l'hypothèse nulle est généralement rejetée lorsque $p \leq 0,05$.

$$H_0 : X \perp\!\!\!\perp Y, \text{ contre } H_1 : X \not\perp\!\!\!\perp Y$$

et,

$$T = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Règle de décision :

- Si $p \leq \alpha$: les variables présentent une association statistiquement significative donc on rejette H_0 .
- Si $p > \alpha$: on ne possède pas suffisamment de preuves pour conclure que les variables sont associées, donc on ne pas rejeter H_0 .

```
> resu.chi2 = chisq.test(president12[1:96,3:12])
> resu.chi2

Pearson's Chi-squared test

data:  president12[1:96, 3:12]
X-squared = 1226242, df = 855, p-value < 2.2e-16
```

D'après le test de dépendance, la p-valeur est très faible. On rejette donc l'hypothèse nulle d'indépendance des variables. Il existe bien un lien entre ces variables.

5. **Effectuer l'analyse factorielle des correspondances (AFC) du tableau avec les choix suivants :**
- lignes actives = les 96 départements de la France métropolitaine
 - colonnes actives = les 10 candidats

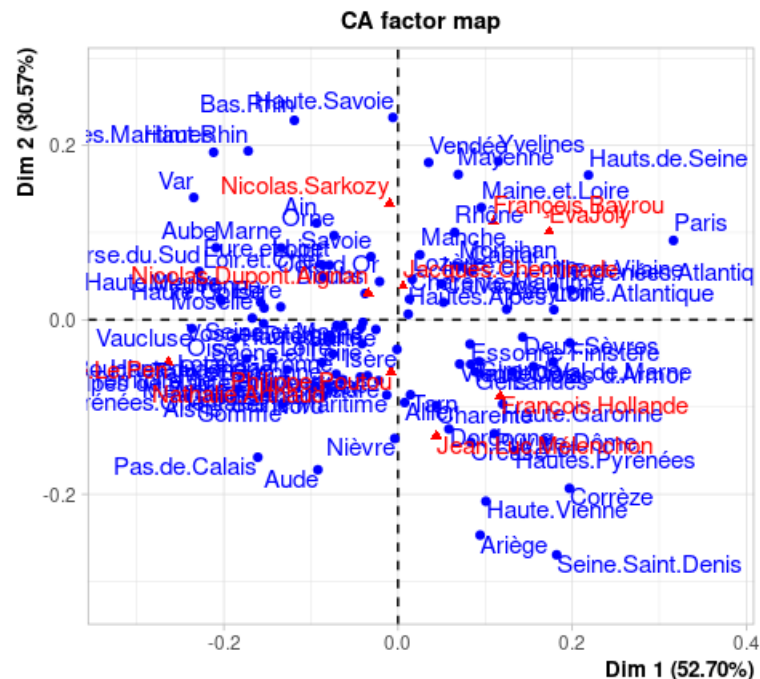
```
1 library(FactoMineR)
2 AFC = CA(president12[1:96, 3:12], row.sup=NULL, col.sup=NULL)
3 AFC
```

```

**Results of the Correspondence Analysis (CA)**
The row variable has 96 categories; the column variable has 10 categories
The chi square of independence between the two variables is equal to 1226242 (p-value = 0 )
*The results are available in the following objects:

  name                description
1  "Seig"              "eigenvalues"
2  "$col"              "results for the columns"
3  "$col$scoord"      "coord. for the columns"
4  "$col$cos2"        "cos2 for the columns"
5  "$col$contrib"     "contributions of the columns"
6  "$row"             "results for the rows"
7  "$row$scoord"      "coord. for the rows"
8  "$row$cos2"        "cos2 for the rows"
9  "$row$contrib"     "contributions of the rows"
10 "$call"             "summary called parameters"
11 "$call$marge.col"  "weights of the columns"
12 "$call$marge.row"  "weights of the rows"

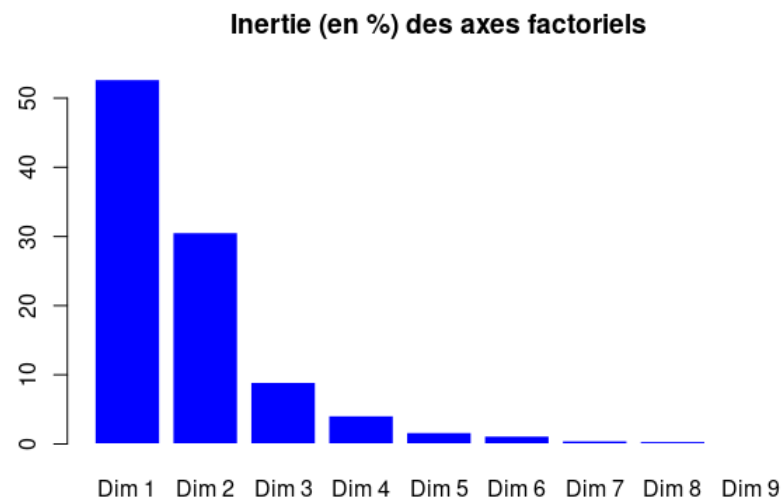
```



Ce graphe permet de visualiser les similitudes de vote entre les différents départements.

6. Créer un diagramme en bâtons pour étudier la décroissance de l'inertie des axes :

```
1 barplot(AFC$eig[,2], names=paste("Dim",1:length(AFC$eig[,2])),
2       main="Inertie (en %) des axes factoriels",
3       col="blue", border="white")
4
```



7. Retrouver l'inertie totale de deux facons :

- a partir du test du 2;
- en calculant la somme des inerties de tous les axes factoriels issus de l'AFC.
- En deduire que la valeur du V de Cramer est egale a 0.063.

```

1  as.numeric(resu.chi2$statistic) / sum(president12[1:96,3:12])
2
3  eigenValue = get_eigenvalue(AFC)
4
5  inertieTotale = as.numeric(sum(eigenValue[,1]))
6  inertieTotale
7
8  VCramer = sqrt(inertieTotale / (min(nrow(president12[1:96,3:12]), ncol(
9  president12[1:96,3:12]))-1))
10 VCramer

```

```

[1] 0.03548686
[1] 0.03548686
[1] 0.06279319

```

8. Calculer la valeur propre moyenne. En deduire le nombre d'axes dont l'inertie est superieure a l'inertie moyenne d'un axe.

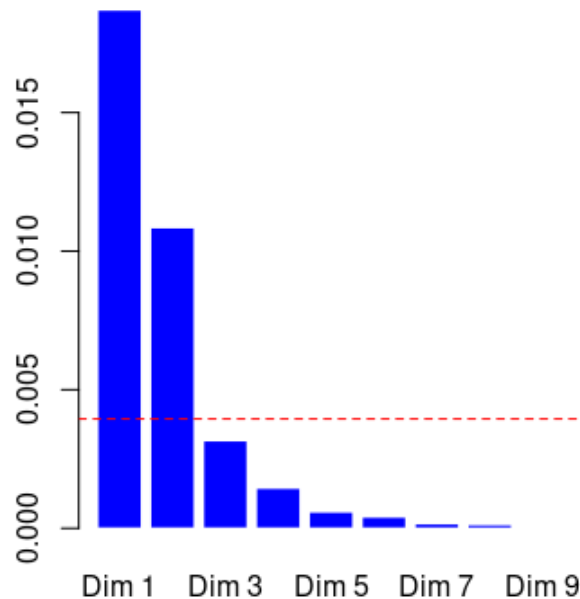
```

1  sum(AFC$eig[,1])/9
2
3  barplot (AFC$eig[1:9,1],
4           names=paste("Dim",1:9),
5           main="Inertie (brute) des axes factoriels",
6           col="blue", border="white", cex.main=1.5)
7  abline(h=sum(AFC$eig[,1])/9, lty=2, col="red")
8  text(6,0.005,"Inertie moyenne",
9       col="red", cex=1.5)

```

```
> sum(AFC$eig[,1])/9 # Valeur propre moyenne  
[1] 0.003942984
```

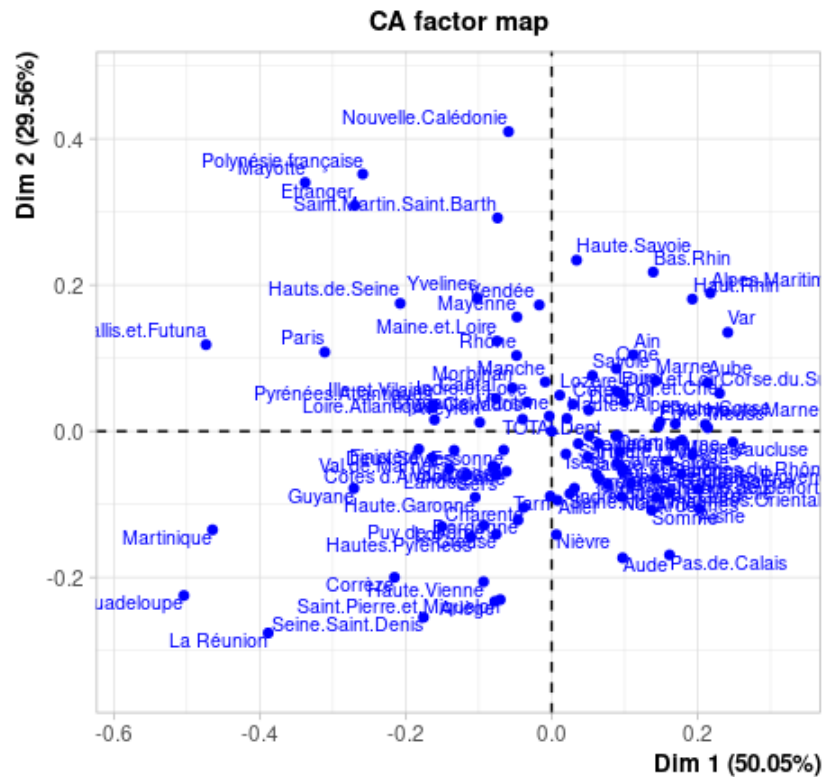
Inertie (brute) des axes factoriels



On peut en conclure qu'il suffit de 2 axes pour que l'inertie soit supérieure à l'inertie moyenne.

9. **Produire le premier plan factoriel (axes 1 et 2) des seuls individus (departements). Remarque : pour rendre le graphique plus lisible, il est possible de reduire la taille des libelles des departements en soumettant a nouveau les commandes de l'AFC apres avoir ajoute l'argument `cex=0.7` a la fin de la commande `plot.CA`.**

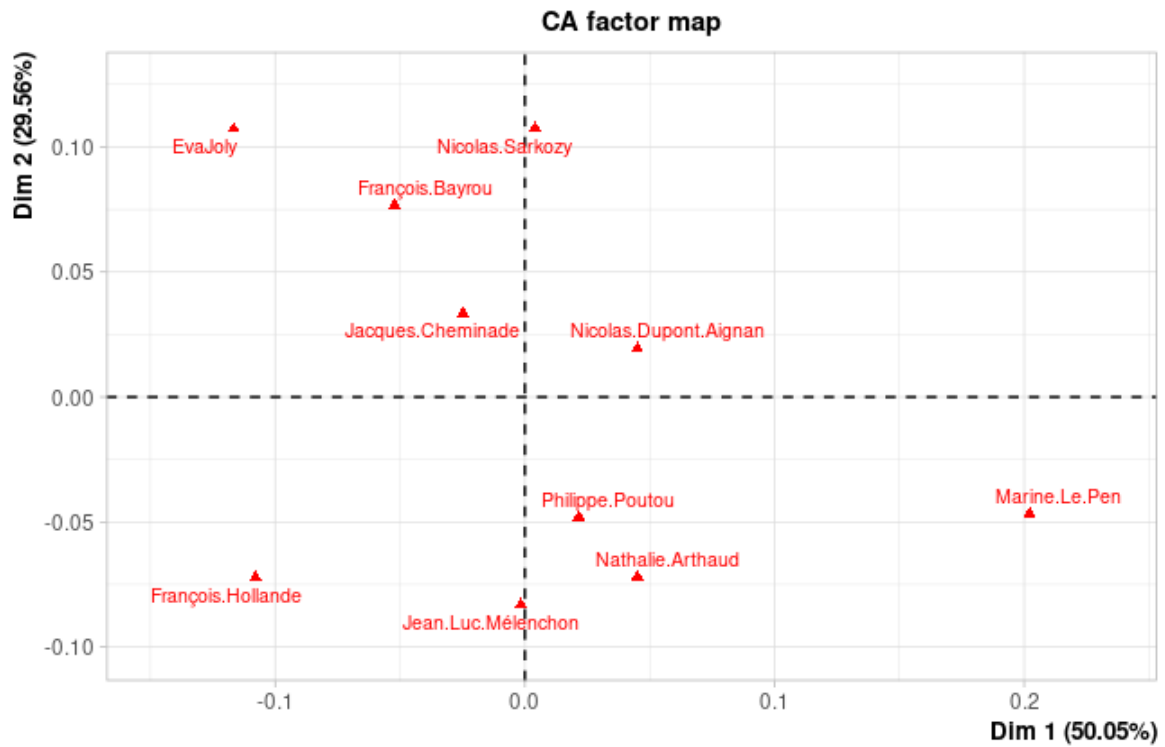
```
1 AFC = CA(president12[1:108, 3:12], graph = FALSE)  
2 plot.CA(AFC, cex=0.7, invisible = "col")
```



L'ACF permet de visualiser les similitudes de vote entre les différents départements.

10. **Produire le premier plan factoriel des colonnes actives (candidats) et supplémentaires.**

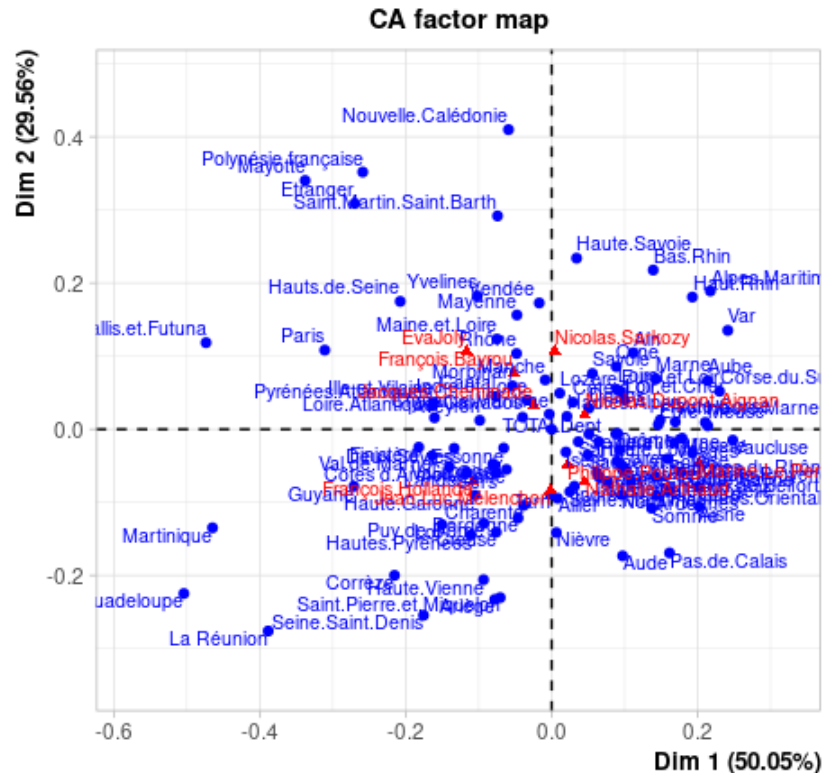
```
1 plot.CA(AFC, cex=0.7, invisible = "row")
```

Le graphe des colonnes nous permet de voir les similitudes entre les candidats. Deux candidats proches sur le graphique ont un profil de votes identiques dans les départements.

11. **Produire le premier plan factoriel avec la totalité des lignes et colonnes (actives et supplémentaires).**

```
1 plot.CA(AFC, cex=0.7)
```



12. Quels sont les départements ayant le plus fortement contribué à la construction de l'axe 1 ? de l'axe 2 ? Pour obtenir les contributions triées par ordre décroissant des axes 1 et 2, on peut soumettre les commandes :

```
1 sort(AFC$row$contrib[,1], T)
2 sort(AFC$row$contrib[,2], T)
```

```
> sort(AFC$row$contrib[,1], T)
      Paris      Hauts.de.Seine      Var      Alpes.Maritimes
1.535816e+01  5.428419e+00  5.115035e+00  4.045516e+00
Loire.Atlantique  Pas.de.Calais  Bouches.du.Rhône  Finistère
3.834126e+00  3.407218e+00  3.371757e+00  3.371119e+00
Val.de.Marne  Ile.et.Vilaine  Seine.Saint.Denis  Vaucluse
2.876997e+00  2.869037e+00  2.739564e+00  2.715024e+00
```

```
> sort(AFC$row$contrib[,2], T)
      Seine.Saint.Denis      Bas.Rhin      Yvelines      Haute.Savoie
1.033938e+01  8.168503e+00  6.337401e+00  5.766694e+00
Alpes.Maritimes  Pas.de.Calais  Hauts.de.Seine  Haut.Rhin
5.731288e+00  5.628983e+00  5.381849e+00  4.049472e+00
Vendée  Var  Nord  Haute.Vienne
3.531328e+00  3.147859e+00  2.789018e+00  2.516891e+00
```

Les départements ayant le plus contribué à l'axe 1 sont :

Paris, Hauts-de-Seine, Var, les Alpes Maritimes et la Loire Atlantique.

Les départements ayant le plus contribué à l'axe 2 sont :

Seine-Saint-Denis, Bas-Rhin, Yvelines, la Haute-Savoie et les Alpes Maritimes.

13. Faire de meme pour les candidats.

```
1 sort(AFC$col$contrib[,1], T)
2 sort(AFC$col$contrib[,2], T)
```

```
> sort(AFC$col$contrib[,1], T)
Marine.Le.Pen      François.Hollande      EvaJoly      François.Bayrou Nicolas.Dupont.Aignan      Nathalie.Arthaud      Philippe.Poutou
64.902393775      29.580049222      2.786568303      2.202574998      0.323948133      0.101554809      0.047596475
Nicolas.Sarkozy    Jacques.Cheminade    Jean.Luc.Mélenchon
0.039199958      0.013526517      0.002587809
> sort(AFC$col$contrib[,2], T)
Nicolas.Sarkozy    François.Hollande    Jean.Luc.Mélenchon    François.Bayrou    Marine.Le.Pen    EvaJoly    Nathalie.Arthaud
47.24443889      22.33617626      11.50434065      8.09484714      5.85284461      3.98969413      0.43787765
Philippe.Poutou    Nicolas.Dupont.Aignan    Jacques.Cheminade
0.39605324      0.10152503      0.04220239
```

Les candidats qui ont le plus contribué à l'axe 1 sont :

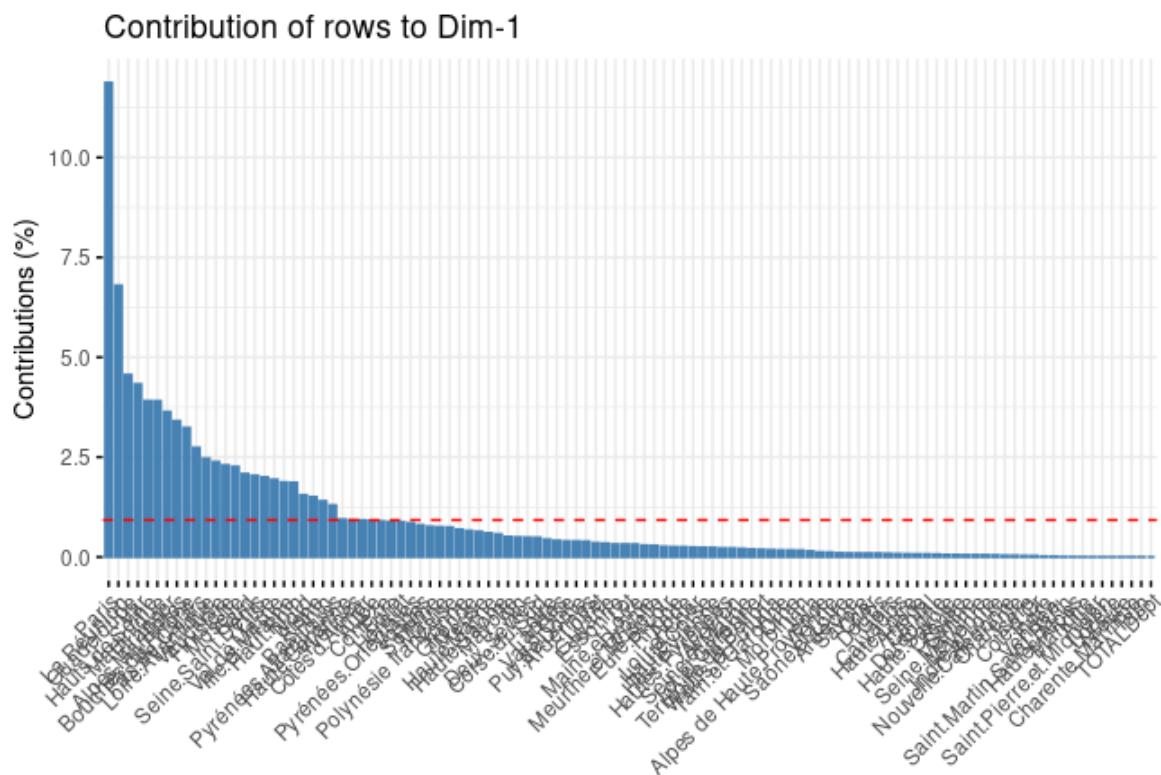
Marine Le Pen et François Hollande.

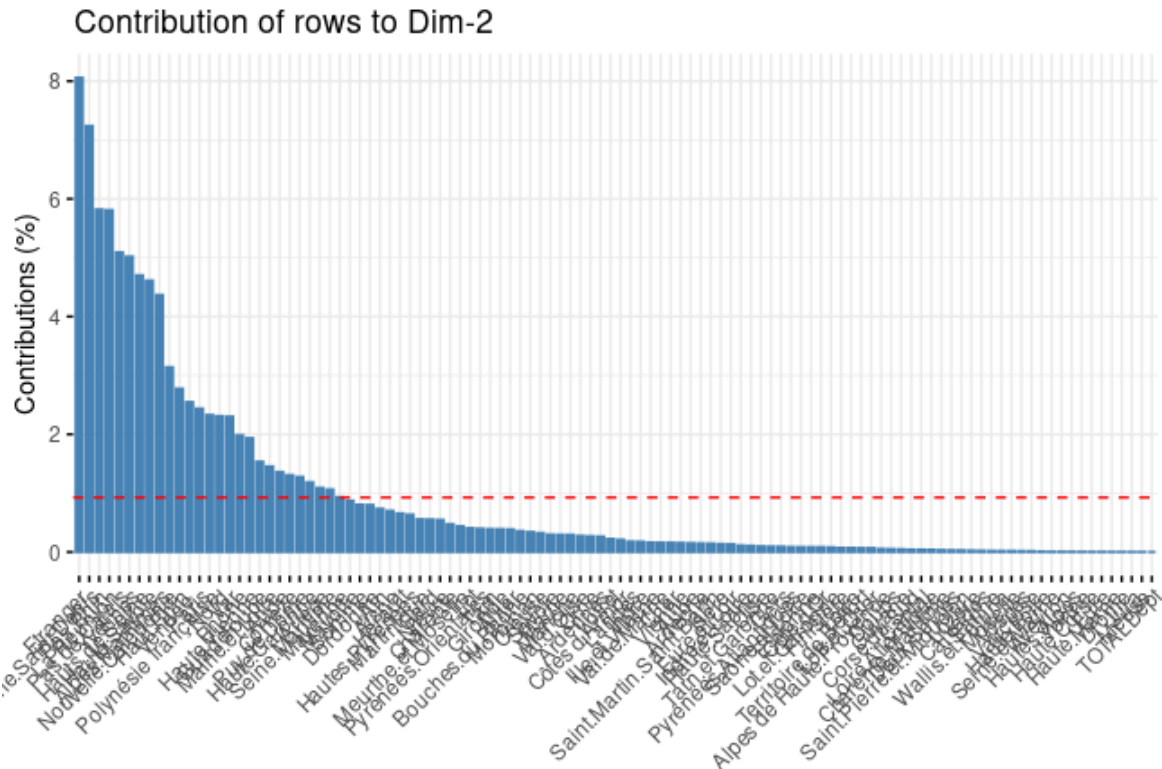
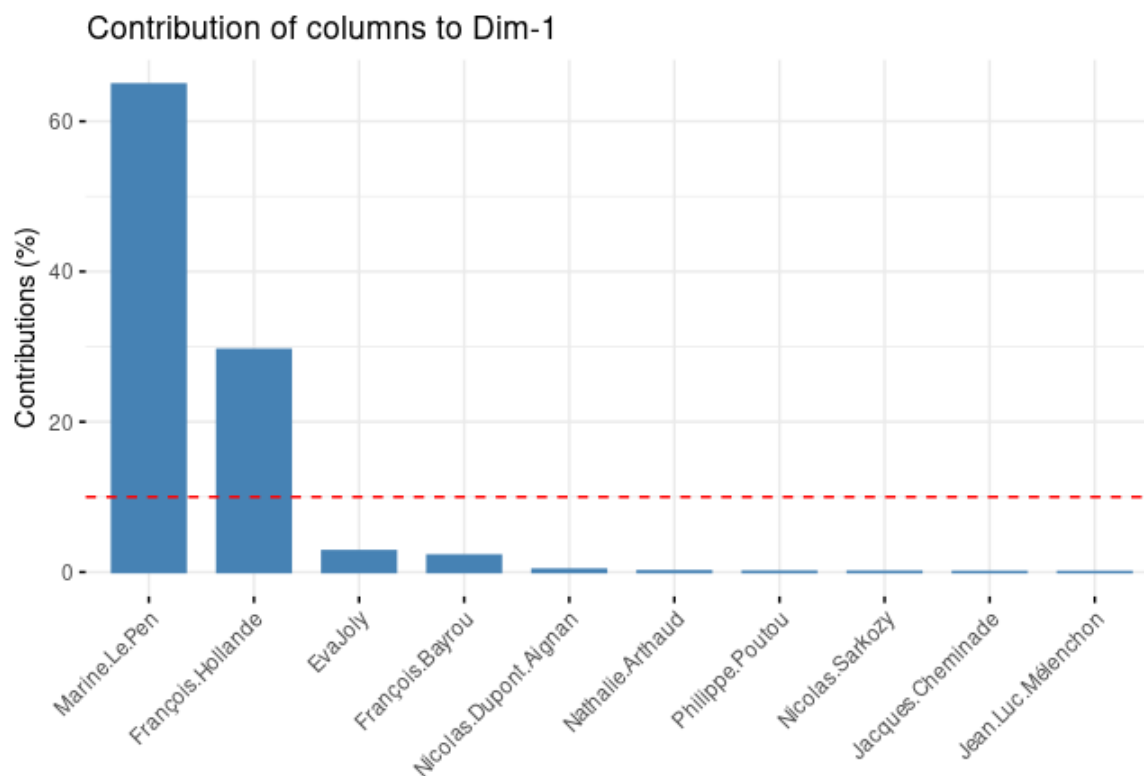
Les candidats ayant le plus contribué à l'axe 2 sont :

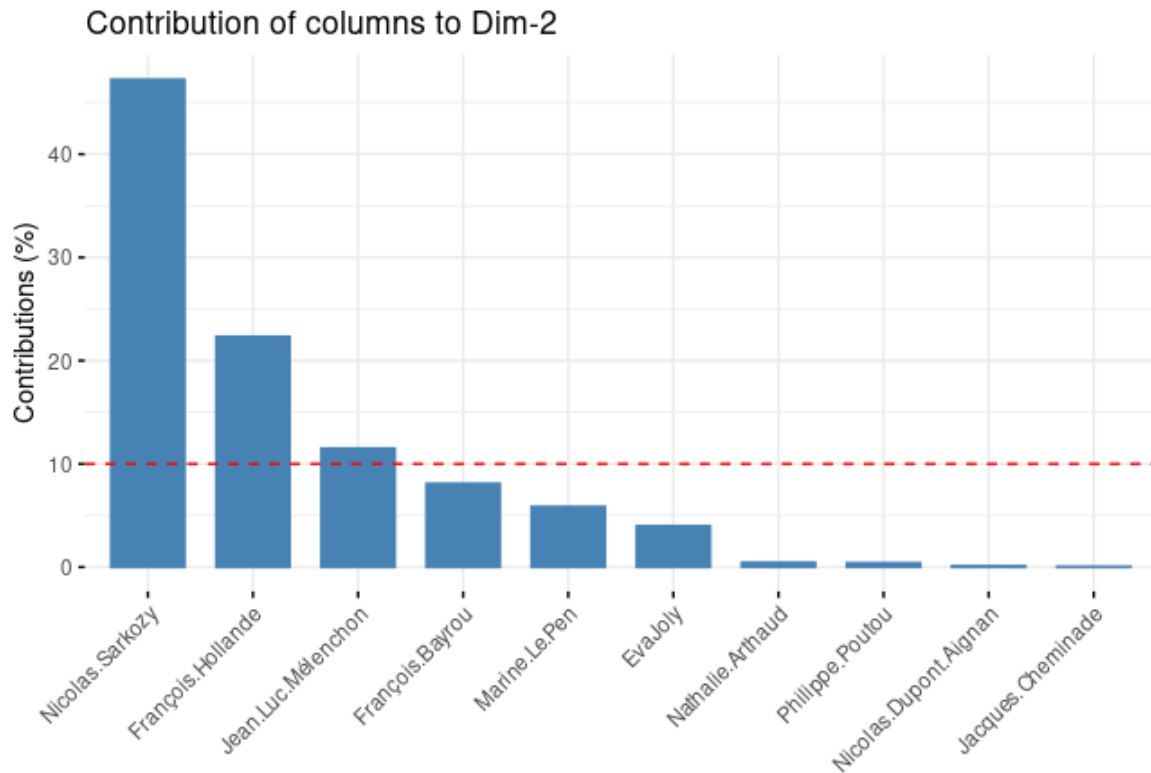
Nicolas Sarkozy et François Hollande.

14. Quels sont les départements et les candidats ayant le plus fortement contribué à l'inertie globale ?

```
1 fviz_contrib(AFC, choice="row", axe=1)
2 fviz_contrib(AFC, choice="col", axe=1)
3 fviz_contrib(AFC, choice="row", axe=2)
4 fviz_contrib(AFC, choice="col", axe=2)
```







Les candidats qui ont le plus contribué sont :

Marine Le Pen, Nicolas Sarkozy et François Hollande.

Les départements ayant le plus contribué sont :

Paris, les Hauts de Seine, la Seine-Saint-Denis et le Var.

15. **Calculer la distance de chacun des départements au centre de gravité du nuage. Comment interpréter la distance la plus courte ?**

```

1 distance = 0
2 for(i in 1:nrow(AFC$row$coord)){
3   distance[i]= sqrt(AFC$row$coord[i,1]^2 + AFC$row$coord[i,2]^2)#Calcul des
  distances
4 }
5 names = row.names(AFC$row$coord)
6
7 df = data.frame(Department = names, Distance_center = distance)
8 df
9
10 df[df["Distance_center"] == min(df["Distance_center"])]

```

Department <chr>	Distance_center <dbl>
Ain	0.15309298
Aisne	0.22805920
Allier	0.09482413
Alpes de Haute.Provence	0.11012843
Hautes.Alpes	0.02730670
Alpes.Maritimes	0.28816967
Ardèche	0.08533303
Ardennes	0.16919743
Ariège	0.24557028
Aube	0.22374941

1-10 of 108 rows

Previous 1 2 3

```
[1] "TOTALDept" "0.00000000"
```

Les Hautes-Alpes sont le département le plus proche au centre de gravité du nuage de points, c'est-à-dire que c'est le département le plus proche du comportement moyen de tous les départements.

16. Interpréter les axes factoriels 1 et 2.

Les départements ayant le plus contribué à l'axe 1 sont : Paris, Hauts-de-Seine, Var, les Alpes-Maritimes et la Loire-Atlantique.

Les candidats qui ont le plus contribué à l'axe 1 sont : Marine Le Pen et François Hollande.

Les départements Paris et Hauts-de-Seine ont une contribution plus importante au pôle positif de la première dimension, tandis que les départements Var et Alpes-Maritimes ont une contribution majeure au pôle négatif de la première dimension. De la même manière, c'est l'opposition entre Marine Le Pen et François Hollande qui semble définir la dimension 1 pour les candidats.

Les départements ayant le plus contribué à l'axe 2 sont : Seine-Saint-Denis, Bas-Rhin, Yvelines, la Haute-Savoie et les Alpes-Maritimes.

Les candidats ayant le plus contribué à l'axe 2 sont : Nicolas Sarkozy et François Hollande

Pareillement, l'axe factoriel 2 semblent être définies par l'opposition des départements Seine-Saint-Denis et du couple Bas Rhin-Yvelines. Côté candidats, c'est l'opposition du couple Nicolas Sarkozy avec François Hollande.

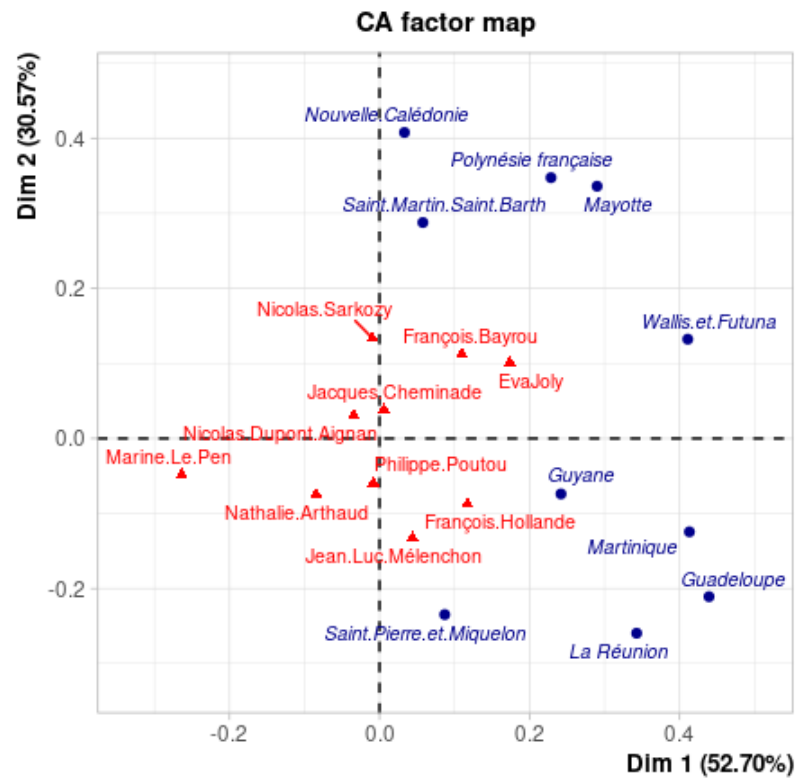
17. Refaire l'AFC avec les départements et territoires d'outremer en elements illustratifs (supplementaires). Interpréter la position de ces derniers sur le plan factoriel (valider votre commentaire en construisant les graphiques des profils-lignes de ces nouveaux départements.).

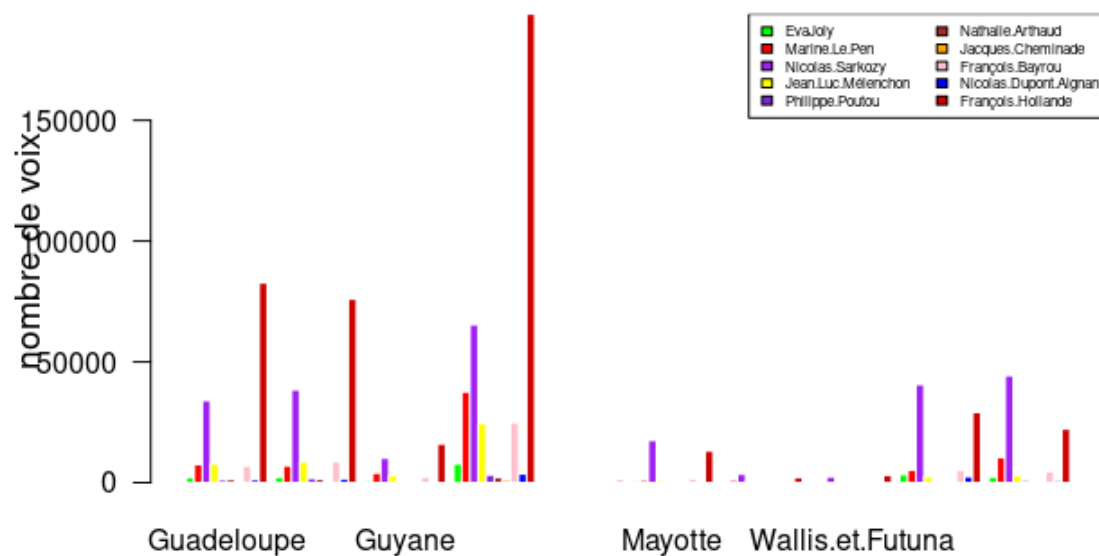
```
1 plot.CA(CA(president12[1:106,3:12], row.sup = 97:106), cex = 0.7, invisible =
  "row")
2
3 color = c("green", "red", "purple", "yellow",
```

```

4      "purple3", "brown", "orange", "pink", "blue", "red3")
5
6  barplot(t(president12[97:106,3:12]),
7          col = color, border="white", horiz=F, beside=T,
8          ylab="nombre de voix", las=1, cex.lab=1.2)
9
10 legend("topright", cex = 0.5, ncol = 2,
11        legend = colnames(president12)[3:12],
12        fill = color)

```





Les départements des outre-mer ont un vote très éloigné de la moyenne globale. La Guyane a en grande partie voté pour Marine Le Pen, tandis que les autres départements ont le plus voté pour Nicolas Sarkozy et François Hollande.